

# Matrix Concentration & Computational Linear Algebra

Short course at École Normale Supérieure, Paris, July 2019

**Joel A. Tropp**

Steele Family Professor of Applied & Computational Mathematics

California Institute of Technology



Typeset on February 23, 2021

Copyright ©2019 Joel A. Tropp

**Cite as:**

Joel A. Tropp, *Matrix Concentration & Computational Linear Algebra*, Caltech CMS Lecture Notes 2019-01, Pasadena, July 2019.

**Available from**

<http://resolver.caltech.edu/CaltechAUTHORS:20190715-125341188>

These lecture notes are composed using an adaptation of a template designed by Mathias Legrand, licensed under CC BY-NC-SA 3.0 (<http://creativecommons.org/licenses/by-nc-sa/3.0/>).



Public domain

# Contents

<b>Preface</b>	<b>ix</b>
<b>Notation</b>	<b>xiii</b>
<b>1</b>	<b>1</b>
<b>1.1</b>	<b>1</b>
<b>1.1.1</b>	<b>1</b>
<b>1.1.2</b>	<b>3</b>
<b>1.2</b>	<b>3</b>
<b>1.2.1</b>	<b>3</b>
<b>1.2.2</b>	<b>4</b>
<b>1.2.3</b>	<b>5</b>
<b>1.2.4</b>	<b>6</b>
<b>1.3</b>	<b>6</b>
<b>1.3.1</b>	<b>7</b>
<b>1.3.2</b>	<b>7</b>
<b>1.4</b>	<b>7</b>
<b>1.4.1</b>	<b>7</b>
<b>1.4.2</b>	<b>8</b>
<b>1.5</b>	<b>9</b>
<b>1.5.1</b>	<b>10</b>
<b>1.5.2</b>	<b>10</b>

<b>1.6</b>	<b>The rectangular case</b>	<b>12</b>
1.6.1	The self-adjoint dilation	12
1.6.2	Rectangular matrix Bernstein	12
<b>2</b>	<b>Matrix Approximation by Sampling</b>	<b>15</b>
<b>2.1</b>	<b>Matrix sampling estimators</b>	<b>15</b>
2.1.1	An error estimate	16
2.1.2	Discussion	18
<b>2.2</b>	<b>Application: Random features</b>	<b>20</b>
2.2.1	Kernel matrices	20
2.2.2	Random features and low-rank approximation of the kernel matrix	21
2.2.3	Examples of random feature maps	23
2.2.4	Error bound for the random feature approximation	24
2.2.5	Analysis of the random feature approximation	25
<b>3</b>	<b>Quantum State Tomography</b>	<b>27</b>
<b>3.1</b>	<b>Postulates of quantum mechanics</b>	<b>27</b>
3.1.1	Recapitulation: Discrete probability theory	27
3.1.2	Noncommutative probability theory	29
3.1.3	Aside: Geometric intuition and the Bloch ball	30
<b>3.2</b>	<b>Quantum state tomography</b>	<b>32</b>
3.2.1	Geometric aspects and measurement design	33
3.2.2	Statistical aspects and convergence	35
<b>3.3</b>	<b>Quantum state tomography via matrix sampling</b>	<b>36</b>
3.3.1	Estimating the bias of a coin	36
3.3.2	The matrix sampling estimator	37
3.3.3	Sample complexity of the sample average	37
3.3.4	Projection onto the set of quantum states	39
3.3.5	Generalization: Projected least squares	41
<b>4</b>	<b>Graph Laplacians</b>	<b>43</b>
<b>4.1</b>	<b>Multigraph basics</b>	<b>43</b>
4.1.1	Undirected multigraphs	43
4.1.2	Connected components	44
4.1.3	Multidegree and total weight	45
4.1.4	Interpretation: Plumbing	45
4.1.5	Interpretation: Resistor networks	45
4.1.6	Example: A random walk	45
<b>4.2</b>	<b>Laplacian basics</b>	<b>46</b>
4.2.1	The Laplacian of a multigraph	46
4.2.2	Correspondence between multigraphs and Laplacians	47
4.2.3	Projectors and pseudoinverses	47
4.2.4	The Dirichlet form	48
4.2.5	Example: Laplacians and cuts	48

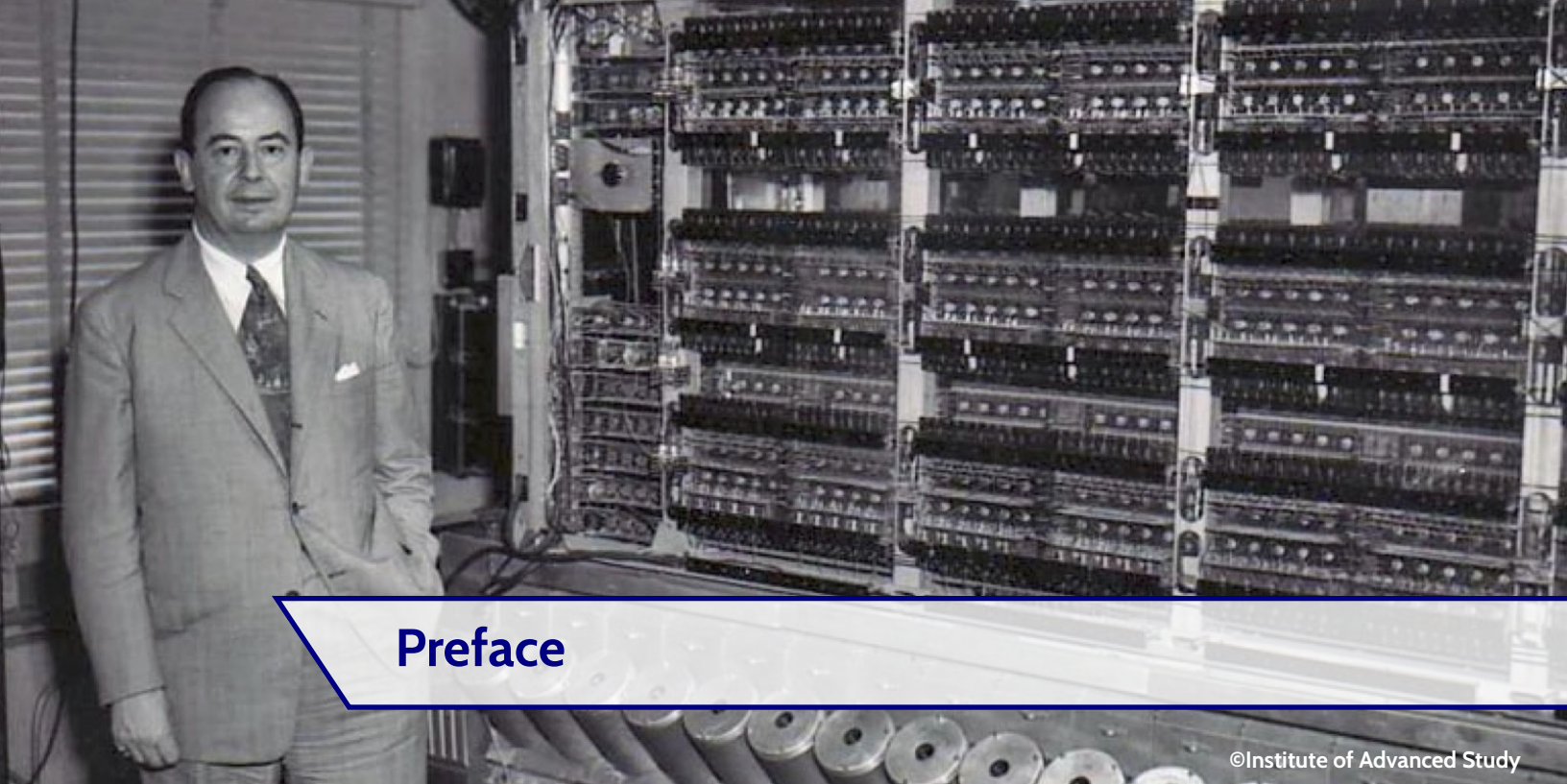
<b>4.3</b>	<b>Harmonic analysis on multigraphs</b>	<b>49</b>
4.3.1	Harmonic functions	49
4.3.2	Example: Hitting probabilities	49
4.3.3	The maximum principle	49
4.3.4	Poles	50
4.3.5	Harmonic extensions	50
4.3.6	Interpretation: Plumbing	51
4.3.7	Interpretation: Resistor networks	52
<b>5</b>	<b>Effective Resistance</b>	<b>53</b>
<b>5.1</b>	<b>Resistance distance</b>	<b>53</b>
5.1.1	Effective resistance	53
5.1.2	Effective resistance is a metric	54
5.1.3	An alternative representation	55
5.1.4	Leverage of a multiedge	56
<b>5.2</b>	<b>Approximating a Laplacian by sampling</b>	<b>57</b>
5.2.1	Spectral approximation	57
5.2.2	The sampling model	58
5.2.3	The sampling probabilities	58
5.2.4	The analysis	58
5.2.5	Computational aspects	59
5.2.6	Conclusion	59
<b>6</b>	<b>Solving Laplacian Systems</b>	<b>61</b>
<b>6.1</b>	<b>Cholesky meets Laplace</b>	<b>61</b>
6.1.1	Setup	61
6.1.2	Laplacian systems	61
6.1.3	Solution via Cholesky decomposition	62
<b>6.2</b>	<b>Cholesky decomposition: Matrix view</b>	<b>62</b>
6.2.1	Setup	62
6.2.2	First step of the Cholesky decomposition	62
6.2.3	Cholesky decomposition, without pivoting	63
6.2.4	Cholesky decomposition, with pivoting	64
6.2.5	Computational cost	64
<b>6.3</b>	<b>Cholesky decomposition: Graph view</b>	<b>64</b>
6.3.1	Setup	64
6.3.2	First step of the Cholesky decomposition	64
6.3.3	Stars and cliques	66
6.3.4	Cholesky decomposition of a Laplacian	67
6.3.5	An opportunity	67
<b>7</b>	<b>Matrix Martingales</b>	<b>69</b>
<b>7.1</b>	<b>Matrix-valued random processes</b>	<b>69</b>
7.1.1	Martingales	69

7.1.2	Matrix martingales . . . . .	70
7.1.3	Adapted sequences . . . . .	70
7.1.4	Stopped processes . . . . .	70
<b>7.2</b>	<b>Tail bounds for matrix-valued processes</b>	<b>70</b>
7.2.1	Corrector processes . . . . .	71
7.2.2	Lower bounds for the supermartingale . . . . .	71
7.2.3	A tail bound for matrix martingales . . . . .	71
<b>7.3</b>	<b>Building a corrector process</b>	<b>72</b>
7.3.1	Correctors . . . . .	73
7.3.2	Lieb's theorem and Tropp's corollary . . . . .	73
7.3.3	Example: The Bernstein corrector . . . . .	73
7.3.4	Example: The Chernoff corrector . . . . .	74
7.3.5	From correctors to corrector processes . . . . .	74
7.3.6	Correctors tensorize . . . . .	74
7.3.7	The composition rule . . . . .	75
<b>7.4</b>	<b>Example: The matrix Freedman inequality</b>	<b>75</b>
<b>8</b>	<b>Sparse Cholesky . . . . .</b>	<b>77</b>
<b>8.1</b>	<b>Approximate solutions of Laplacian systems</b>	<b>77</b>
8.1.1	Approximate solutions . . . . .	77
8.1.2	Approximate Cholesky decomposition . . . . .	78
8.1.3	Preconditioning . . . . .	78
8.1.4	Summary . . . . .	78
<b>8.2</b>	<b>Overview of the algorithm</b>	<b>79</b>
8.2.1	Setup . . . . .	79
8.2.2	The SparseCholesky procedure . . . . .	79
8.2.3	Laplacian approximations . . . . .	81
<b>8.3</b>	<b>Preliminaries for the analysis</b>	<b>81</b>
8.3.1	The normalizing map . . . . .	81
8.3.2	The approximation requirement . . . . .	82
8.3.3	Splitting the edges . . . . .	82
<b>8.4</b>	<b>Sampling from a clique</b>	<b>83</b>
8.4.1	Setup . . . . .	83
8.4.2	Eliminating a vertex . . . . .	84
8.4.3	The sampling procedure . . . . .	84
8.4.4	Expectation of the random multiedge . . . . .	85
8.4.5	Each multiedge has bounded leverage . . . . .	85
8.4.6	Corrector for the random multiedge . . . . .	86
8.4.7	An unbiased estimator for the clique . . . . .	87
8.4.8	The clique induced by a random vertex . . . . .	87
8.4.9	Corrector for the clique estimator . . . . .	88

<b>8.5</b>	<b>Analysis of SparseCholesky</b>	<b>89</b>
8.5.1	A stopping time . . . . .	89
8.5.2	The approximate Schur complements . . . . .	89
8.5.3	The corrector process . . . . .	90
8.5.4	The martingale tail bound . . . . .	90
8.5.5	The running time . . . . .	91
8.5.6	The grand finale . . . . .	91
	<b>Further Reading</b> . . . . .	<b>93</b>
	<b>Bibliography</b> . . . . .	<b>97</b>







## Preface

©Institute of Advanced Study

Over the last decade, random matrices have become ubiquitous in applied and computational mathematics. As this trend accelerates, more and more researchers must confront random matrices as part of their work. Classical random matrix theory can be difficult to use, and it is often silent about the questions that come up in modern applications. As a consequence, it has become imperative to develop new tools that are easy to use and that apply to a wide range of random matrices.

### Matrix concentration inequalities

Matrix concentration inequalities are among the most popular of these new methods. For a self-adjoint random matrix  $\mathbf{Y}$  with expectation  $\mathbb{E} \mathbf{Y}$ , matrix concentration theorems provide probabilistic bounds on quantities like

$$\|\mathbf{Y} - \mathbb{E} \mathbf{Y}\|.$$

The symbol  $\|\cdot\|$  always refers the spectral norm, also known as the  $\ell_2$  operator norm. Bounds of this form give us a lot of information about how the random matrix  $\mathbf{Y}$  is related to its expectation  $\mathbb{E} \mathbf{Y}$ . In particular,

- Each linear functional of  $\mathbf{Y}$  is close to the same linear functional of  $\mathbb{E} \mathbf{Y}$ .
- Each eigenvalue of  $\mathbf{Y}$  is close to the corresponding eigenvalue of  $\mathbb{E} \mathbf{Y}$ .
- Each eigenvector of  $\mathbf{Y}$  is close to the corresponding eigenvector of  $\mathbb{E} \mathbf{Y}$  when the eigenvalue is isolated from the rest of the spectrum.
- We can bound the expected norm of the random matrix:

$$\|\mathbf{Y}\| = \|\mathbb{E} \mathbf{Y}\| \pm \|\mathbf{Y} - \mathbb{E} \mathbf{Y}\|.$$

The last point is, perhaps, the most interesting. Indeed, norm bounds for random matrices are quite valuable by themselves, and they used to be rather hard to obtain before the matrix concentration technology was developed.

Matrix concentration results for self-adjoint random matrices also have formal consequences for rectangular random matrices. We will focus on the self-adjoint case because it is more fundamental, and it already supports many fascinating applications.

### Random matrix models

Without additional information about the random matrix  $Y$ , we cannot hope to say anything interesting. This work treats two basic, but very fruitful, models for the random matrix.

First, the *independent sum model* posits that

$$Y = \sum_{i=1}^n X_i \quad \text{where } \{X_i\} \text{ is statistically independent.}$$

This model captures a wide range of examples. The most classical is the sample covariance matrix; see [Tro15, Chap. 1] for discussion and analysis. In this course, we will explore more modern examples from machine learning, quantum information theory, and combinatorics.

Second, we will consider the *matrix martingale model*, where

$$Y_k = \sum_{i=1}^k X_i \quad \text{is a martingale.}$$

This model offers a powerful lens for studying the behavior of iterative randomized algorithms in linear algebra. The main purpose of this course is to show how concentration for matrix martingales supported the development and analysis of an efficient algorithm for solving graph Laplacian linear systems.

### Other applications of matrix concentration

Matrix concentration tools have already found a place in many areas of the mathematical sciences, including

- numerical linear algebra [Tro11b]
- numerical analysis [MB17]
- uncertainty quantification [CG14]
- statistics [Kol11]
- econometrics [CC13]
- approximation theory [CDL13]
- sampling theory [BG13]
- machine learning [DKC13; Lop+14]
- learning theory [FSV12; MKR12]
- math signal processing [Che+14]
- optimization [CSW12]
- graphics and vision [HCG14]
- quantum information [Hol12]
- algorithms [HO14; Kyn17]
- combinatorics [Oli10]
- *et cetera*.

These references are chosen more or less at random from a long menu of possibilities. See the monograph [Tro15] for an overview of the main results on matrix concentration, many detailed applications, and additional background references. Other recommendations for further reading appear at the end of these notes.

## About this course

These lecture notes were written to support the short course

### *Matrix Concentration & Computational Linear Algebra*

delivered by the author at École Normale Supérieure in Paris from 1–5 July 2019 as part of the summer school “High-dimensional probability and algorithms.”

The aim of this course is to present some practical computational applications of matrix concentration. Lecture 1 provides a brief treatment of the matrix Bernstein inequality, which is the most valuable single result about matrix concentration. We apply this result to study several empirical matrix approximations:

- Random feature approximation of a kernel matrix (Lecture 2).
- Linear estimators for quantum state tomography (Lecture 3).
- Sparse approximation of a combinatorial graph (Lecture 5).

Our primary goal is to develop a complete treatment of a near-linear time algorithm for solving a linear system in a graph Laplacian matrix. This remarkable algorithm was developed by Rasmus Kyng and Sushant Sachdeva [KS16], following earlier work [Kyn+16] by Dan Spielman’s group. The algorithm closely resembles the classic incomplete Cholesky decomposition, and I believe that it is likely to have an impact on computational practice. Our presentation of this result takes place in steps:

- Harmonic analysis on graphs (Lecture 4).
- Interpretation of graphs as resistor networks (Lecture 5).
- Cholesky factorization of a graph Laplacian (Lecture 6).
- Theory of matrix martingales (Lecture 7).
- The SparseCholesky algorithm (Lecture 8).

In my opinion, the SparseCholesky algorithm is the most spectacular application of matrix concentration. I doubt that it could have been developed before the foundations of matrix concentration were in place.

## Prerequisites

Since the audience of this short course consists primarily of French graduate students and researchers, I have assumed a moderate level of mathematical and computational preparation:

- Intermediate linear algebra [Axl15], including experience with positive-semidefinite matrices and the semidefinite order [Bha97; Bha07].
- Elementary numerical linear algebra [TB97], including Cholesky decomposition, solution of triangular systems, and the conjugate gradient algorithm.
- Intermediate probability, including basic scalar concentration inequalities [BLM13] and real-valued discrete-time martingales [Wil91].
- Elementary spectral graph theory [Spi12].

We will develop the background for most of the applications in sufficient detail that no additional preparation is needed.

**Caveat lector**

These notes do not meet the standard of a scholarly publication. Here are some issues that you should be aware of.

- Some of the text has been copied and pasted directly from my own published work (Lectures 1, 2, and 7).
- The notes for Lecture 3 were written primarily by Richard Kueng on the basis of our joint work [Guh+18]. I take responsibility for any mistakes that appear.
- The treatment of graph theory is my interpretation of Dan Spielman's lecture notes [Spi] and Rasmus Kyng's dissertation [Kyn17]. Any errors are mine.
- Owing to the varied provenance of the material, the mathematical notation may not be fully consistent among different lectures.
- I have tried to provide citations for the main results, but these notes are largely devoid of references, historical background, and context.
- These notes have only received a cursory proofreading.

**Why is there a photo of von Neumann?**

You may be wondering why John von Neumann greets you at the door of this Preface. In 1947, von Neumann and Goldstine [NG47] developed the foundations for rounding error analysis. They formulated Gaussian elimination and Cholesky decomposition as triangular matrix factorizations. They showed how to analyze the numerical properties of the linear system solver based on this approach. In a 1951 follow-up paper [GN51], they proposed a random matrix model for the rounding errors in these computations. See [Grc11] for a gloss on this research.

The landmark papers of von Neumann and Goldstine are among the earliest works on solving linear systems on a computer, and they are the first to bring random matrix theory in contact with computational linear algebra. I cannot think of a more suitable genie to inhabit these notes.

**Acknowledgments**

The summer school “High-dimensional probability and algorithms” was funded by Université PSL and CNRS. I would like to thank the organizers, Claire Boyer, Djalil Chafaï, and Joseph Lehec, for an engaging week. Additional funding for my research and this course was provided by ONR Awards N00014-17-12146 and N00014-18-12363.

The computational cost of the **SparseCholesky** algorithm was miscalculated in the original manuscript because of stray parentheses; the costs are slightly higher than reported. This version corrects the error, which was reported by Rasmus Kyng.

Joel A. Tropp

[jtropp@cms.caltech.edu](mailto:jtropp@cms.caltech.edu)

<http://users.cms.caltech.edu/~jtropp>

Steele Family Professor of Applied & Computational Mathematics  
California Institute of Technology  
Pasadena, California  
July 2019



## Notation

I have selected notation that is common in the linear algebra and probability literature. I have tried to be consistent in using the symbols that are presented below. There are some minor variations in different lectures, including the letter that indicates the dimension of a matrix and the indexing of sums.

### Linear algebra

We work in a real or complex linear space. The letters  $d$  and  $n$  (and occasionally others) are used to denote the dimension of this space, which is always finite. For example, we write  $\mathbb{R}^d$  or  $\mathbb{C}^n$ . Matrix concentration results apply equally in the real and complex setting. We may write  $\mathbb{F}$  to refer to either field, or we may omit the field entirely.

We use the delta notation for standard basis vectors:  $\delta_i$  has a one in the  $i$ th coordinate and zeros elsewhere. The vector  $\mathbf{1}$  has ones in each entry. The dimension of these vectors is determined by context.

The symbol  $*$  denotes the (conjugate) transpose of a vector or a matrix. We equip  $\mathbb{F}^d$  with the standard inner product  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^* \mathbf{y}$ . The inner product generates the Euclidean norm  $\|\mathbf{x}\|^2 = \langle \mathbf{x}, \mathbf{x} \rangle$ .

We write  $\mathbb{H}_d(\mathbb{F})$  for the real-linear space of  $d \times d$  self-adjoint matrices with entries in the field  $\mathbb{F}$ . Recall that a matrix is self-adjoint when  $\mathbf{A} = \mathbf{A}^*$ . We equip the space  $\mathbb{H}_d$  with the trace inner product  $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{tr}(\mathbf{X}\mathbf{Y})$ , which generates the Frobenius norm  $\|\mathbf{X}\|_{\mathbb{F}}^2 = \langle \mathbf{X}, \mathbf{X} \rangle$ . The map  $\text{tr}[\cdot]$  returns the trace of a square matrix; we instate the convention that nonlinear functions bind before the trace.

A self-adjoint matrix with dimension  $d$  has  $d$  real eigenvalues, with an associated orthonormal set of  $d$  eigenvectors. The maps  $\lambda_{\min}(\cdot)$  and  $\lambda_{\max}(\cdot)$  return the minimum and maximum eigenvalues of a self-adjoint matrix. The symbol  $\mathbf{I}$  denotes the identity matrix; its dimensions are determined by context.

A self-adjoint matrix is *positive semidefinite* (psd) if its eigenvalues are nonnegative; a self-adjoint matrix is *positive definite* (pd) if its eigenvalues are positive. The symbol  $\preceq$  refers to the psd order:  $\mathbf{A} \preceq \mathbf{H}$  if and only if  $\mathbf{H} - \mathbf{A}$  is psd.

We define a *standard matrix function* on a self-adjoint matrix using the eigenvalue decomposition. For any  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,

$$\mathbf{A} = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^* \quad \text{implies} \quad f(\mathbf{A}) = \sum_{i=1}^n f(\lambda_i) \mathbf{u}_i \mathbf{u}_i^*.$$

When we apply a real function to a self-adjoint matrix, we are always referring to the associated standard matrix function. In particular, we often encounter powers, exponentials, and logarithms.

Occasionally, we need the linear space  $\mathbb{M}^{d_1 \times d_2}(\mathbb{F})$  of  $d_1 \times d_2$  matrices over the field  $\mathbb{F}$ . In this context, the symbol  $\|\cdot\|$  also refers to the  $\ell_2$  operator norm.

We write  $\text{lin}$  for the linear hull of a family of vectors. The operators  $\text{range}$  and  $\text{null}$  extract the range and null space of a matrix. The operator  $^\dagger$  extracts the pseudoinverse.

### Probability

The map  $\mathbb{P}\{\cdot\}$  returns the probability of an event. The operator  $\mathbb{E}[\cdot]$  returns the expectation of a random variable taking values in a linear space. We only include the brackets when it is necessary for clarity, and we impose the convention that nonlinear functions bind before the expectation.

### Graphs

A multigraph  $\mathbf{G}$  has a ground set  $\mathbf{V}$  of  $n$  vertices. A multiedge is an undirected pair  $e = uv = \{u, v\}$  of vertices. A multigraph involves a set  $\mathbf{E}$  of  $m$  multiedges, which may involve many edges connecting the same pair of vertices. The absolute value  $|\cdot|$  returns the cardinality of a set of vertices or a set of edges.

We write  $\mathbb{R}^{\mathbf{V}}$  for the set of real-valued functions on the set  $\mathbf{V}$  of vertices. The symbol  $\mathbb{H}_{\mathbf{V}}$  refers to the linear space of (real) self-adjoint matrices acting on  $\mathbb{R}^{\mathbf{V}}$ . We may identify these linear spaces with  $\mathbb{R}^n$  and  $\mathbb{H}_n(\mathbb{R})$ .

The notation  $u \sim v$  means that two vertices are neighbors. The notations  $u \in e$  and  $e \ni u$  both mean that the multiedge  $e$  is incident on (i.e., contains) the vertex  $u$ .

The degree  $\deg(u, \mathbf{G})$  of a vertex  $u$  in a multigraph  $\mathbf{G}$  is the total number of multiedges incident on  $u$ . The total weight  $w_{\mathbf{G}}(u)$  of a vertex  $u$  is the sum of the weights of the multiedges incident on  $u$ .

We reserve the letter  $\mathbf{L}$  for the Laplacian matrix of the multigraph  $\mathbf{G}$ . The symbol  $\Phi$  denotes the normalizing map associated with this Laplacian:

$$\Phi(\mathbf{M}) = (\mathbf{L}^\dagger)^{1/2} \mathbf{M} (\mathbf{L}^\dagger)^{1/2}.$$

The exponent  $^{1/2}$  extracts the unique psd square root of a psd matrix. The number  $\varrho(u, v)$  is the effective resistance between vertices  $u$  and  $v$ .

### Order notation

We use the familiar order notation from computer science. The symbol  $\Theta(\cdot)$  refers to asymptotic equality. The symbol  $O(\cdot)$  refers to an asymptotic upper bound.

# 1. Matrix Concentration

©1999–2003 by Jamie Zawinski

Most of the text in this lecture is copied from my monograph [Tro15, Chaps. 3, 5, 6].

This lecture contains the analysis that delivers exponential matrix concentration inequalities. The approach that we take can be viewed as a matrix extension of the Laplace transform method, sometimes referred to as the “Bernstein trick.” In the scalar setting, this trick (*soi disant*) is one of the most basic and successful paths to reach concentration inequalities for sums of independent random variables. It turns out that there is a very satisfactory version of this argument that applies to sums of independent random matrices. In the general setting, however, we must invest more care and wield sharper tools to execute this technique.

## 1.1 The matrix Laplace transform method

In the scalar setting, the Laplace transform method allows us to obtain tail bounds for a random variable in terms of its mgf. The starting point for our theory is the observation that a similar result holds in the matrix setting.

### 1.1.1 Tail bounds

First, we introduce the Laplace transform method for bounding the extreme eigenvalues of a self-adjoint matrix. This approach to matrix concentration was proposed by Ahlswede & Winter [AW02]. We present a formulation and proof developed by Roberto Oliveira [Oli10].

**Proposition 1.1 (Tail bounds for eigenvalues).** Let  $Y$  be a random self-adjoint matrix. For

all  $t \in \mathbb{R}$ ,

$$\mathbb{P} \{ \lambda_{\max}(\mathbf{Y}) \geq t \} \leq \inf_{\theta > 0} e^{-\theta t} \mathbb{E} \operatorname{tr} e^{\theta \mathbf{Y}}; \quad (1.1)$$

$$\mathbb{P} \{ \lambda_{\min}(\mathbf{Y}) \leq t \} \leq \inf_{\theta < 0} e^{-\theta t} \mathbb{E} \operatorname{tr} e^{\theta \mathbf{Y}}. \quad (1.2)$$

In words, we can control the tail probabilities of the extreme eigenvalues of a random matrix by producing a bound for the *trace* of the matrix mgf. The proof of this fact parallels the classical argument, but there is a twist.

*Proof.* We begin with (1.1). Fix a positive number  $\theta$ , and observe that

$$\mathbb{P} \{ \lambda_{\max}(\mathbf{Y}) \geq t \} = \mathbb{P} \left\{ e^{\theta \lambda_{\max}(\mathbf{Y})} \geq e^{\theta t} \right\} \leq e^{-\theta t} \mathbb{E} e^{\theta \lambda_{\max}(\mathbf{Y})}.$$

The first identity holds because  $a \mapsto e^{\theta a}$  is a monotone increasing function, so the event does not change under the mapping. The second relation is Markov's inequality. To control the exponential, note that

$$e^{\theta \lambda_{\max}(\mathbf{Y})} = e^{\lambda_{\max}(\theta \mathbf{Y})} = \lambda_{\max}(e^{\theta \mathbf{Y}}) \leq \operatorname{tr} e^{\theta \mathbf{Y}}. \quad (1.3)$$

The first identity holds because the maximum eigenvalue is a positive-homogeneous map. The second depends on the spectral mapping theorem. The inequality follows because the exponential of an self-adjoint matrix is positive definite. The maximum eigenvalue of a positive-definite matrix is dominated by the trace. Combine the latter two displays to reach

$$\mathbb{P} \{ \lambda_{\max}(\mathbf{Y}) \geq t \} \leq e^{-\theta t} \mathbb{E} \operatorname{tr} e^{\theta \mathbf{Y}}.$$

This inequality is valid for any positive  $\theta$ , so we may take an infimum to achieve the tightest possible bound.

To prove (1.2), we use a similar approach. Fix a negative number  $\theta$ , and calculate that

$$\mathbb{P} \{ \lambda_{\min}(\mathbf{Y}) \leq t \} = \mathbb{P} \left\{ e^{\theta \lambda_{\min}(\mathbf{Y})} \geq e^{\theta t} \right\} \leq e^{-\theta t} \mathbb{E} e^{\theta \lambda_{\min}(\mathbf{Y})} = e^{-\theta t} \mathbb{E} e^{\lambda_{\max}(\theta \mathbf{Y})}.$$

The function  $a \mapsto e^{\theta a}$  reverses the inequality in the event because it is monotone decreasing. The last identity depends on the relationship between minimum and maximum eigenvalues. Finally, we introduce the inequality (1.3) for the trace exponential and minimize over negative values of  $\theta$ . ■

In the proof of Proposition 1.1, it may seem crude to bound the maximum eigenvalue by the trace. It turns out that, at most, this estimate results in a loss of a factor that is logarithmic in the dimension of the matrix. At the same time, our maneuver allows us to exploit some amazing convexity properties of the trace exponential.



### 1.1.2 Expectation bounds

We can adapt the proof of Proposition 1.1 to obtain bounds for the expectation of the maximum eigenvalue of a random self-adjoint matrix. This argument is somewhat less interesting in the scalar setting, where it states that the exponential mean of a random variable is an upper bound for the arithmetic mean.

**Proposition 1.2 (Expectation bounds for eigenvalues).** Let  $Y$  be a random self-adjoint matrix. Then

$$\mathbb{E} \lambda_{\max}(Y) \leq \inf_{\theta > 0} \frac{1}{\theta} \log \mathbb{E} \operatorname{tr} e^{\theta Y}; \quad (1.4)$$

$$\mathbb{E} \lambda_{\min}(Y) \geq \sup_{\theta < 0} \frac{1}{\theta} \log \mathbb{E} \operatorname{tr} e^{\theta Y}. \quad (1.5)$$

*Proof.* We establish the bound (1.4); the proof of (1.5) is quite similar. Fix a positive number  $\theta$ , and calculate that

$$\begin{aligned} \mathbb{E} \lambda_{\max}(Y) &= \frac{1}{\theta} \mathbb{E} \log e^{\lambda_{\max}(\theta Y)} \leq \frac{1}{\theta} \log \mathbb{E} e^{\lambda_{\max}(\theta Y)} \\ &= \frac{1}{\theta} \log \mathbb{E} \lambda_{\max}(e^{\theta Y}) \leq \frac{1}{\theta} \log \mathbb{E} \operatorname{tr} e^{\theta Y}. \end{aligned}$$

The first identity holds because the maximum eigenvalue is a positive-homogeneous map. The second relation is Jensen's inequality. The third follows when we use the spectral mapping theorem to draw the exponential inside the eigenvalue map. The final inequality depends on the fact that the trace of a positive-definite matrix dominates the maximum eigenvalue. ■

## 1.2 Matrix moments and cumulants

At the heart of the Laplace transform method are the moment generating function (mgf) and the cumulant generating function (cgf) of a random variable. In this section, we define these functions rigorously, and we explore some of their properties.

### 1.2.1 The matrix mgf and cgf

We begin by presenting matrix versions of the mgf and cgf.

**Definition 1.3 (Matrix mgf and cgf).** Let  $X$  be a random self-adjoint matrix. The *matrix moment generating function*  $M_X$  and the *matrix cumulant generating function*  $\Xi_X$  are given by

$$M_X(\theta) = \mathbb{E} e^{\theta X} \quad \text{and} \quad \Xi_X(\theta) = \log \mathbb{E} e^{\theta X} \quad \text{for } \theta \in \mathbb{R}. \quad (1.6)$$

Note that the expectations may not exist for all values of  $\theta$ .

The matrix mgf  $M_X$  and matrix cgf  $\Xi_X$  contain information about the distribution of the random matrix  $X$ , including its mean and variance. Propositions 1.1 and 1.2 show how to exploit the data encoded in these functions to control the eigenvalues.

Let us dilate on Definition 1.3. Observe that the matrix mgf and cgf have formal power series expansions:

$$\mathbf{M}_X(\theta) = \mathbf{I} + \sum_{q=1}^{\infty} \frac{\theta^q}{q!} (\mathbb{E} \mathbf{X}^q) \quad \text{and} \quad \mathbf{\Xi}_X(\theta) = \sum_{q=1}^{\infty} \frac{\theta^q}{q!} \mathbf{\Psi}_q.$$

We call the coefficients  $\mathbb{E} \mathbf{X}^q$  *matrix moments*, and we refer to  $\mathbf{\Psi}_q$  as a *matrix cumulant*. The matrix cumulant  $\mathbf{\Psi}_q$  has a formal expression as a (noncommutative) polynomial in the matrix moments up to order  $q$ . In particular, the first cumulant is the mean and the second cumulant is the variance:

$$\mathbf{\Psi}_1 = \mathbb{E} \mathbf{X} \quad \text{and} \quad \mathbf{\Psi}_2 = \mathbb{E} \mathbf{X}^2 - (\mathbb{E} \mathbf{X})^2.$$

Higher-order matrix cumulants are harder to write down and interpret.

### 1.2.2 The failure of the matrix mgf

We would like to use the Laplace transform bounds from Section 1.1 to study a sum of independent random matrices. In the scalar setting, the Laplace transform method is effective for studying independent sums because the mgf and the cgf decompose. In the matrix case, the situation is more subtle, and the goal of this section is to indicate where things go awry.

Consider an independent sequence  $\{X_k\}$  of real random variables. The mgf of the sum satisfies a multiplication rule:

$$M_{(\sum_k X_k)}(\theta) = \mathbb{E} \exp\left(\sum_k \theta X_k\right) = \mathbb{E} \prod_k e^{\theta X_k} = \prod_k \mathbb{E} e^{\theta X_k} = \prod_k M_{X_k}(\theta). \quad (1.7)$$

The first identity is the definition of an mgf. The second relation holds because the exponential map converts a sum of real scalars to a product, and the third relation requires the independence of random variables. The last identity, again, is the definition.

At first, we might imagine that a similar relationship holds for the matrix mgf. Consider an independent sequence  $\{\mathbf{X}_k\}$  of random self-adjoint matrices. Perhaps,

$$\mathbf{M}_{(\sum_k \mathbf{X}_k)}(\theta) \stackrel{?}{=} \prod_k \mathbf{M}_{\mathbf{X}_k}(\theta). \quad (1.8)$$

Unfortunately, this hope shatters when we subject it to interrogation.

It is not hard to find the reason that (1.8) fails. The identity (1.7) depends on the fact that the scalar exponential converts a sum into a product. In contrast, for self-adjoint matrices,

$$e^{\mathbf{A}+\mathbf{H}} \neq e^{\mathbf{A}}e^{\mathbf{H}} \quad \text{unless } \mathbf{A} \text{ and } \mathbf{H} \text{ commute.}$$

If we introduce the trace, the situation improves somewhat:

$$\text{tr } e^{\mathbf{A}+\mathbf{H}} \leq \text{tr } e^{\mathbf{A}}e^{\mathbf{H}} \quad \text{for all self-adjoint } \mathbf{A}, \mathbf{H}. \quad (1.9)$$

The result (1.9) is known as the Golden–Thompson inequality, a famous theorem from statistical physics. Unfortunately, the analogous bound may fail for three matrices:

$$\text{tr } e^{\mathbf{A}+\mathbf{H}+\mathbf{T}} \not\leq \text{tr } e^{\mathbf{A}}e^{\mathbf{H}}e^{\mathbf{T}} \quad \text{for certain self-adjoint } \mathbf{A}, \mathbf{H}, \mathbf{T}.$$

It seems that we have reached an impasse.

What if we consider the cgf instead? The cgf of a sum of independent real random variables satisfies an addition rule:

$$\Xi_{(\sum_k X_k)}(\theta) = \log \mathbb{E} \exp \left( \sum_k \theta X_k \right) = \log \prod_k \mathbb{E} e^{\theta X_k} = \sum_k \Xi_{X_k}(\theta). \quad (1.10)$$

The relation (1.10) follows when we extract the logarithm of the multiplication rule (1.7). This result looks like a more promising candidate for generalization because a sum of self-adjoint matrices remains self-adjoint. We might hope that

$$\Xi_{(\sum_k X_k)}(\theta) \stackrel{?}{=} \sum_k \Xi_{X_k}(\theta).$$

As stated, this putative identity also fails. Nevertheless, the addition rule (1.10) admits a very satisfactory extension to matrices. In contrast to the scalar case, the proof involves much deeper considerations.

### 1.2.3 A theorem of Lieb

To find the appropriate generalization of the addition rule for cgfs, we turn to the literature on matrix analysis. Here, we discover a famous result of Elliott Lieb on the convexity properties of the trace exponential function.

**Theorem 1.4 (Lieb, 1973).** Fix a self-adjoint matrix  $\mathbf{H}$  with dimension  $d$ . The function

$$\mathbf{A} \mapsto \text{tr} \exp(\mathbf{H} + \log \mathbf{A})$$

is a concave map on the convex cone of  $d \times d$  positive-definite matrices.

In the scalar case, the analogous function  $a \mapsto \exp(h + \log a)$  is linear, so this result describes a new type of phenomenon that emerges when we move to the matrix setting. See [Tro15, Chap. 8] for a complete proof of Theorem 1.4 from first principles.

Lieb's theorem is valuable to us because the Laplace transform bounds from Section 1.1 involve the trace exponential function. To highlight the connection, let us rephrase Theorem 1.4 in probabilistic terms.

**Corollary 1.5 (Tropp, 2010).** Let  $\mathbf{H}$  be a fixed self-adjoint matrix, and let  $\mathbf{X}$  be a random self-adjoint matrix of the same dimension. Then

$$\mathbb{E} \text{tr} \exp(\mathbf{H} + \mathbf{X}) \leq \text{tr} \exp(\mathbf{H} + \log \mathbb{E} e^{\mathbf{X}}).$$

*Proof.* Introduce the random matrix  $\mathbf{Y} = e^{\mathbf{X}}$ . Then

$$\begin{aligned} \mathbb{E} \text{tr} \exp(\mathbf{H} + \mathbf{X}) &= \mathbb{E} \text{tr} \exp(\mathbf{H} + \log(\mathbf{Y})) \\ &\leq \text{tr} \exp(\mathbf{H} + \log(\mathbb{E} \mathbf{Y})) = \text{tr} \exp(\mathbf{H} + \log \mathbb{E} e^{\mathbf{X}}). \end{aligned}$$

The first identity follows from the interpretation of the matrix logarithm as the functional inverse of the matrix exponential for positive-definite matrices. Theorem 1.4 shows that the trace function is concave in  $\mathbf{Y}$ , so Jensen's inequality allows us to draw the expectation inside the function. ■

### 1.2.4 Subadditivity of the matrix cgf

We are now prepared to generalize the addition rule (1.10) for scalar cgfs to the matrix setting. The following result is fundamental to our approach to random matrices.

**Lemma 1.6 (Subadditivity of matrix cgfs).** Consider a finite sequence  $\{X_k\}$  of independent, random, self-adjoint matrices of the same dimension. Then

$$\mathbb{E} \operatorname{tr} \exp \left( \sum_k \theta X_k \right) \leq \operatorname{tr} \exp \left( \sum_k \log \mathbb{E} e^{\theta X_k} \right) \quad \text{for } \theta \in \mathbb{R}. \quad (1.11)$$

Equivalently,

$$\operatorname{tr} \exp \left( \Xi_{(\sum_k X_k)}(\theta) \right) \leq \operatorname{tr} \exp \left( \sum_k \Xi_{X_k}(\theta) \right) \quad \text{for } \theta \in \mathbb{R}. \quad (1.12)$$

The parallel between the additivity rule (1.10) and the subadditivity rule (1.12) is striking. With our level of preparation, it is easy to prove this result. We just apply the bound from Corollary 1.5 repeatedly.

*Proof.* Without loss of generality, we assume that  $\theta = 1$  by absorbing the parameter into the random matrices. Let  $\mathbb{E}_k$  denote the expectation with respect to  $X_k$ , the remaining random matrices held fixed. Abbreviate

$$\Xi_k = \log \mathbb{E}_k e^{X_k} = \log \mathbb{E} e^{X_k}.$$

We may calculate that

$$\begin{aligned} \mathbb{E} \operatorname{tr} \exp \left( \sum_{k=1}^n X_k \right) &= \mathbb{E} \mathbb{E}_n \operatorname{tr} \exp \left( \sum_{k=1}^{n-1} X_k + X_n \right) \\ &\leq \mathbb{E} \operatorname{tr} \exp \left( \sum_{k=1}^{n-1} X_k + \log (\mathbb{E}_n e^{X_n}) \right) \\ &= \mathbb{E} \mathbb{E}_{n-1} \operatorname{tr} \exp \left( \sum_{k=1}^{n-2} X_k + X_{n-1} + \Xi_n \right) \\ &\leq \mathbb{E} \mathbb{E}_{n-2} \operatorname{tr} \exp \left( \sum_{k=1}^{n-2} X_k + \Xi_{n-1} + \Xi_n \right) \\ &\dots \leq \operatorname{tr} \exp \left( \sum_{k=1}^n \Xi_k \right). \end{aligned}$$

We use the statistical independence of  $\{X_i\}$  to introduce the iterated expectation. At each step  $m = 1, 2, 3, \dots, n$ , we invoke Corollary 1.5 with the fixed matrix  $H$  equal to

$$H_m = \sum_{k=1}^{m-1} X_k + \sum_{k=m+1}^n \Xi_k.$$

This argument is legitimate because  $H_m$  is independent from  $X_m$ .

The formulation (1.12) follows from (1.11) when we substitute the expression (1.6) for the matrix cgf and make some algebraic simplifications. ■

## 1.3 Master bounds for sums of independent random matrices

We are now prepared to present some general results on the behavior of a sum of independent random matrices. In the next section, we derive some concrete matrix concentration inequalities using this approach.

### 1.3.1 The master inequalities

To obtain the main abstract results, we simply combine the Laplace transform bounds with the subadditivity of the matrix cgf.

**Theorem 1.7 (Master bounds for a sum of independent random matrices).** Consider a finite sequence  $\{\mathbf{X}_k\}$  of independent, random, self-adjoint matrices of the same size. Then

$$\mathbb{E} \lambda_{\max} \left( \sum_k \mathbf{X}_k \right) \leq \inf_{\theta > 0} \frac{1}{\theta} \log \operatorname{tr} \exp \left( \sum_k \log \mathbb{E} e^{\theta \mathbf{X}_k} \right); \quad (1.13)$$

$$\mathbb{E} \lambda_{\min} \left( \sum_k \mathbf{X}_k \right) \geq \sup_{\theta < 0} \frac{1}{\theta} \log \operatorname{tr} \exp \left( \sum_k \log \mathbb{E} e^{\theta \mathbf{X}_k} \right). \quad (1.14)$$

Furthermore, for all  $t \in \mathbb{R}$ ,

$$\mathbb{P} \left\{ \lambda_{\max} \left( \sum_k \mathbf{X}_k \right) \geq t \right\} \leq \inf_{\theta > 0} e^{-\theta t} \operatorname{tr} \exp \left( \sum_k \log \mathbb{E} e^{\theta \mathbf{X}_k} \right); \quad (1.15)$$

$$\mathbb{P} \left\{ \lambda_{\min} \left( \sum_k \mathbf{X}_k \right) \leq t \right\} \leq \inf_{\theta < 0} e^{-\theta t} \operatorname{tr} \exp \left( \sum_k \log \mathbb{E} e^{\theta \mathbf{X}_k} \right). \quad (1.16)$$

*Proof.* Substitute the subadditivity rule for matrix cgfs, Lemma 1.6, into the two matrix Laplace transform results, Proposition 1.1 and Proposition 1.2. ■

### 1.3.2 Additional tools

To use Theorem 1.7, we need semidefinite bounds on the matrix cgf that reflect structural properties of the random matrices that appear in the sum. To implement this program, we need several basic facts from matrix analysis.

**Fact 1.8 (Trace exponential is monotone).** If  $\mathbf{A} \preceq \mathbf{H}$ , then  $\operatorname{tr} \exp(\mathbf{A}) \leq \operatorname{tr} \exp(\mathbf{H})$ . ■

**Fact 1.9 (Logarithm is operator monotone).** If  $\mathbf{A} \preceq \mathbf{H}$ , then  $\log \mathbf{A} \preceq \log \mathbf{H}$ . ■

See [Tro15, Chap. 8] for the proofs of these results.

As a consequence of Fact 1.8, it suffices to produce semidefinite upper bounds for the matrix cgfs that appear in the formulas of Theorem 1.7. As a consequence of Fact 1.9, we can obtain a semidefinite upper bound for the matrix cgf from a semidefinite upper bound for the matrix mgf. We will see these ideas in action in the next section.

## 1.4 Example: Matrix Bernstein

We continue with the matrix Bernstein inequality, the matrix concentration result that has found the widest application. This result concerns a sum of independent zero-mean random matrices that are subject to a uniform norm bound.

### 1.4.1 Bernstein cgf bound

The first step in using Theorem 1.7 is to develop an estimate for the cgf of a bounded, zero-mean random matrix. This argument closely follows the analog argument in the scalar setting.

**Lemma 1.10 (Bernstein cgf).** Suppose that  $\mathbf{X}$  is a random self-adjoint matrix that satisfies

$$\mathbb{E} \mathbf{X} = \mathbf{0} \quad \text{and} \quad \|\mathbf{X}\| \leq 1.$$

Then

$$\log \mathbb{E} e^{\theta \mathbf{X}} \leq \frac{\theta^2/2}{1 - |\theta|/3} \cdot \mathbb{E} \mathbf{X}^2.$$

*Proof.* Suppose that  $x \in [-1, +1]$ . Using the Taylor series expansion of the exponential,

$$\begin{aligned} e^{\theta x} &= 1 + \theta x + \sum_{p=2}^{\infty} \frac{\theta^p}{p!} x^p \\ &\leq 1 + \theta x + \left( \sum_{p=2}^{\infty} \frac{|\theta|^p}{2 \cdot 3^{p-2}} \right) \cdot x^2 \\ &= 1 + \theta x + \frac{\theta^2/2}{1 - |\theta|/3} \cdot x^2. \end{aligned}$$

Since each eigenvalue of  $\mathbf{X}$  lies in the interval  $[-1, +1]$ , we can apply this inequality to each eigenvalue of  $\mathbf{X}$  to obtain

$$e^{\theta \mathbf{X}} \leq \mathbf{I} + \theta \mathbf{X} + \frac{\theta^2/2}{1 - |\theta|/3} \cdot \mathbf{X}^2.$$

Take the expectation:

$$\mathbb{E} e^{\theta \mathbf{X}} \leq \mathbf{I} + \frac{\theta^2/2}{1 - |\theta|/3} \cdot \mathbb{E} \mathbf{X}^2.$$

Invoke Fact 1.9:

$$\log \mathbb{E} e^{\theta \mathbf{X}} \leq \log \left( \mathbf{I} + \frac{\theta^2/2}{1 - |\theta|/3} \cdot \mathbb{E} \mathbf{X}^2 \right) \leq \frac{\theta^2/2}{1 - |\theta|/3} \cdot \mathbb{E} \mathbf{X}^2.$$

The last relation follows when we apply the numerical inequality  $\log(1 + x) \leq x$ , valid for  $x > -1$ , to each eigenvalue. ■

### 1.4.2 The matrix Bernstein inequality

Combining the master tail bound, Theorem 1.7, with the cgf bound, Lemma 1.10, we arrive at the matrix Bernstein inequality.

**Theorem 1.11 (Matrix Bernstein).** Consider a statistically independent sequence  $\{\mathbf{X}_k : 1 \leq k \leq n\}$  of random matrices with dimension  $d$ . Suppose that

$$\mathbb{E} \mathbf{X}_k = \mathbf{0} \quad \text{and} \quad \|\mathbf{X}_k\| \leq B \quad \text{for each index } k.$$

Introduce the sum of the random matrices:

$$\mathbf{Y} = \sum_{k=1}^n \mathbf{X}_k.$$

Define the matrix variance proxy:

$$\sigma^2 = \|\mathbb{E} \mathbf{Y}^2\| = \left\| \sum_{k=1}^n \mathbb{E} \mathbf{X}_k^2 \right\|.$$

Then, for all  $t \geq 0$ ,

$$\mathbb{P} \{ \|\mathbf{Y}\| \geq t \} \leq 2d \cdot \exp \left( \frac{-t^2/2}{\sigma^2 + Bt/3} \right).$$

Furthermore,

$$\mathbb{E} \|\mathbf{Y}\| \leq \sqrt{2\sigma^2 \log(2d)} + \frac{1}{3} B \log(2d).$$

*Proof.* First, rescale so that  $B = 1$ . The general form of the result follows from homogeneity arguments. The Bernstein cgf bound, Lemma 1.10, implies that

$$\log \mathbb{E} e^{\theta \mathbf{X}_k} \leq g(\theta)(\mathbb{E} \mathbf{X}_k^2) \quad \text{where} \quad g(\theta) = \frac{\theta^2/2}{1 - |\theta|/3}.$$

Note that  $g(\theta) \geq 0$  for all  $\theta \in \mathbb{R}$ .

Substitute these cgf bounds into the master inequality (1.15) to obtain

$$\begin{aligned} \mathbb{P} \{ \lambda_{\max}(\mathbf{Y}) \geq t \} &\leq \inf_{\theta > 0} e^{-\theta t} \operatorname{tr} \exp \left( g(\theta) \sum_k \mathbb{E} \mathbf{X}_k^2 \right) \\ &\leq d \inf_{\theta > 0} e^{-\theta t} \lambda_{\max} \left( \exp \left( g(\theta)(\mathbb{E} \mathbf{Y}^2) \right) \right) \\ &= d \inf_{\theta > 0} e^{-\theta t} \exp \left( g(\theta) \sigma^2 \right). \end{aligned}$$

The first inequality depends on Fact 1.8. Afterward, we bound the trace by the dimension times the maximum eigenvalue. Next, we invoke the spectral mapping theorem and the fact that  $g(\theta) > 0$  to draw the maximum eigenvalue inside the exponential. Identify the variance proxy  $\sigma^2$  by noting that the maximum eigenvalue of the psd matrix  $\mathbb{E} \mathbf{Y}^2$  coincides with its spectral norm.

Finally, we make the clever choice  $\theta = t/(\sigma^2 + t/3)$  to see that

$$\mathbb{P} \{ \lambda_{\max}(\mathbf{Y}) \geq t \} \leq d \cdot \exp \left( \frac{-t^2/2}{\sigma^2 + t/3} \right).$$

An identical argument yields a corresponding probability bound for the minimum eigenvalue of  $\mathbf{Y}$ . Combine the two results with the union bound to arrive at the stated probability bound for the spectral norm. ■

**Exercise 1.1** Use the master inequalities (1.13) and (1.14) to establish the expectation bound that appears in Theorem 1.11.

## 1.5 Example: Matrix Chernoff

As a second example, we develop bounds for the extreme eigenvalues of an independent sum of bounded, psd matrices.

### 1.5.1 Chernoff cgf bound

The matrix Chernoff inequality is based on the following cgf bound. It is a matrix version of a scalar argument.

**Lemma 1.12 (Chernoff cgf).** Suppose that  $\mathbf{X}$  is a random self-adjoint matrix that satisfies

$$\mathbf{0} \preceq \mathbf{X} \preceq \mathbf{I}.$$

Then

$$\log \mathbb{E} e^{\theta \mathbf{X}} \preceq (\mathrm{e}^\theta - 1)(\mathbb{E} \mathbf{X}) \quad \text{for } \theta \in \mathbb{R}.$$

This result is based on a classic computation for real random variables. The matrix extension first appeared in the proof of [AW02, Thm. 19]. See also [Tro12, Lem. 5.8].

*Proof.* The function  $x \mapsto \mathrm{e}^{\theta x}$  is convex, so the graph lies below the chord connecting two points. In particular,

$$\mathrm{e}^{\theta x} \leq 1 + (\mathrm{e}^\theta - 1)x \quad \text{for } x \in [0, 1].$$

The eigenvalues of  $\mathbf{X}$  lie in the interval  $[0, 1]$ , so

$$\mathrm{e}^{\theta \mathbf{X}} \preceq \mathbf{I} + (\mathrm{e}^\theta - 1)\mathbf{X}.$$

Take the expectation:

$$\mathbb{E} \mathrm{e}^{\theta \mathbf{X}} \preceq \mathbf{I} + (\mathrm{e}^\theta - 1)(\mathbb{E} \mathbf{X}).$$

Arguing as in the proof of Lemma 1.10,

$$\log \mathbb{E} \mathrm{e}^{\theta \mathbf{X}} \preceq (\mathrm{e}^\theta - 1)(\mathbb{E} \mathbf{X}).$$

We have used Fact 1.9 and the numerical inequality  $\log(1 + x) \leq x$ . ■

### 1.5.2 Matrix Chernoff inequalities

Combining the master tail bound, Theorem 1.7, with the cgf bound, Lemma 1.12, we arrive at the matrix Chernoff inequalities.

**Theorem 1.13 (Matrix Chernoff).** Consider a statistically independent sequence  $\{\mathbf{X}_k : 1 \leq k \leq n\}$  of random matrices with dimension  $d$ . Suppose that

$$\mathbf{0} \preceq \mathbf{X}_k \preceq B \mathbf{I} \quad \text{for each index } k.$$

Introduce the sum of the random matrices:

$$\mathbf{Y} = \sum_{k=1}^n \mathbf{X}_k.$$

Define the lower and upper eigenvalues of the expectation:

$$\mu_{\min} = \lambda_{\min}(\mathbb{E} \mathbf{Y}) \quad \text{and} \quad \mu_{\max} = \lambda_{\max}(\mathbb{E} \mathbf{Y}).$$



Then

$$\begin{aligned}\mathbb{P}\{\lambda_{\min}(\mathbf{Y}) \leq (1 - \delta) \mu_{\min}\} &\leq d \cdot \left( \frac{e^{-\delta}}{(1 - \delta)^{1-\delta}} \right)^{\mu_{\min}/B} \quad \text{for } 0 < \delta \leq 1; \\ \mathbb{P}\{\lambda_{\max}(\mathbf{Y}) \geq (1 + \delta) \mu_{\max}\} &\leq d \cdot \left( \frac{e^{\delta}}{(1 + \delta)^{1+\delta}} \right)^{\mu_{\max}/B} \quad \text{for } \delta > 0.\end{aligned}$$

*Proof of Theorem 1.13, maximum eigenvalue bound.* We begin with the tail bound for the maximum eigenvalue  $\lambda_{\max}(\mathbf{Y})$ . By a scaling argument, we may assume that  $B = 1$ . The Chernoff cgf bound, Lemma 1.12, implies that

$$\log \mathbb{E} e^{\theta \mathbf{X}_k} \leq g(\theta)(\mathbb{E} \mathbf{X}_k) \quad \text{where} \quad g(\theta) = e^{\theta} - 1.$$

Note that  $g(\theta) > 0$  for  $\theta > 0$ .

Using Fact 1.8, we substitute these cgf bounds into the master inequality (1.15) to reach

$$\begin{aligned}\mathbb{P}\{\lambda_{\max}(\mathbf{Y}) \geq t\} &\leq \inf_{\theta > 0} e^{-\theta t} \operatorname{tr} \exp \left( g(\theta) \sum_k \mathbb{E} \mathbf{X}_k \right) \\ &\leq \inf_{\theta > 0} e^{-\theta t} d \lambda_{\max} \left( \exp(g(\theta)(\mathbb{E} \mathbf{Y})) \right) \\ &= d \inf_{\theta > 0} e^{-\theta t} \exp \left( g(\theta) \lambda_{\max}(\mathbb{E} \mathbf{Y}) \right) \\ &\leq d \inf_{\theta > 0} e^{-\theta t} \exp \left( g(\theta) \mu_{\max} \right).\end{aligned}$$

In the second line, we use the fact that the matrix exponential is pd to bound the trace by  $d$  times the maximum eigenvalue; we have also identified the sum as  $\mathbb{E} \mathbf{Y}$ . The third line follows from the spectral mapping theorem. Next, we use the fact that the maximum eigenvalue is a positive-homogeneous map, which depends on the observation that  $g(\theta) > 0$  for  $\theta > 0$ . Finally, we identify the statistic  $\mu_{\max}$ .

To complete the proof, make the change of variables  $t \mapsto (1 + \delta) \mu_{\max}$ . Then the infimum is achieved at  $\theta = \log(1 + \delta)$ , which leads to the upper tail bound. ■

The lower bounds follow from a related argument that is slightly more delicate.

*Proof of Theorem 1.13, minimum eigenvalue bound.* We now establish the bound for the minimum eigenvalue  $\lambda_{\min}(\mathbf{Y})$ . As before, rescale so that  $B = 1$ . The Chernoff cgf bound, Lemma 1.12, implies that

$$\log \mathbb{E} e^{\theta \mathbf{X}_k} \leq g(\theta)(\mathbb{E} \mathbf{X}_k) \quad \text{where} \quad g(\theta) = e^{\theta} - 1.$$

Note that  $g(\theta) < 0$  when  $\theta < 0$ .

Introduce these cgf bounds into the master inequality (1.16) to reach

$$\begin{aligned}\mathbb{P}\{\lambda_{\min}(\mathbf{Y}) \leq t\} &\leq \inf_{\theta < 0} e^{-\theta t} \operatorname{tr} \exp \left( g(\theta) \sum_k \mathbb{E} \mathbf{X}_k \right) \\ &\leq \inf_{\theta < 0} e^{-\theta t} d \lambda_{\min} \left( \exp(g(\theta)(\mathbb{E} \mathbf{Y})) \right) \\ &= \inf_{\theta < 0} e^{-\theta t} d \exp \left( g(\theta) \lambda_{\min}(\mathbb{E} \mathbf{Y}) \right) \\ &\leq d \inf_{\theta < 0} e^{-\theta t} \exp \left( g(\theta) \cdot \mu_{\min} \right).\end{aligned}$$

The justifications here are similar to those in the previous argument. The only noteworthy point is that we must replace the maximum eigenvalue map with the minimum eigenvalue map because  $g(\theta) < 0$  for  $\theta < 0$ .

Finally, we make the change of variables  $t \mapsto (1 - \delta) \mu_{\min}$ . The infimum is attained at  $\theta = \log(1 - \delta)$ , which yields the lower tail bound. ■

**Exercise 1.2** Derive the following consequences of Theorem 1.13. For  $\delta \in (0, 1]$ ,

$$\begin{aligned}\mathbb{P} \{ \lambda_{\min}(\mathbf{Y}) \leq (1 - \delta) \mu_{\min} \} &\leq d \cdot e^{-\delta^2 \mu_{\min} / (2B)}, \\ \mathbb{P} \{ \lambda_{\max}(\mathbf{Y}) \geq (1 + \delta) \mu_{\max} \} &\leq d \cdot e^{-\delta^2 \mu_{\max} / (3B)}.\end{aligned}$$

These simplifications are often more tractable in practice.

## 1.6 The rectangular case

In these lectures, we will only be using matrix concentration for self-adjoint matrices. Nevertheless, it is important to be aware that concentration results for rectangular matrices follow as a formal consequence. This section outlines the approach.

### 1.6.1 The self-adjoint dilation

The *self-adjoint dilation*  $\mathcal{H}(\mathbf{S})$  of a rectangular matrix  $\mathbf{S} \in \mathbb{M}^{d_1 \times d_2}$  is the self-adjoint matrix

$$\mathcal{H}(\mathbf{S}) := \begin{bmatrix} \mathbf{0} & \mathbf{S} \\ \mathbf{S}^* & \mathbf{0} \end{bmatrix} \in \mathbb{H}_{d_1+d_2}. \quad (1.17)$$

Note that the map  $\mathcal{H}$  is real-linear. By direct calculation,

$$\mathcal{H}(\mathbf{S})^2 = \begin{bmatrix} \mathbf{S}\mathbf{S}^* & \mathbf{0} \\ \mathbf{0} & \mathbf{S}^*\mathbf{S} \end{bmatrix}. \quad (1.18)$$

We also have the spectral identity

$$\lambda_{\max}(\mathcal{H}(\mathbf{S})) = \|\mathcal{H}(\mathbf{S})\| = \|\mathbf{S}\|. \quad (1.19)$$

This point follows from some linear algebraic considerations.

### 1.6.2 Rectangular matrix Bernstein

Using the device of the self-adjoint dilation, we can develop a version of the matrix Bernstein inequality for rectangular matrices.

**Corollary 1.14 (Rectangular matrix Bernstein).** Consider a statistically independent sequence  $\{\mathbf{S}_k : 1 \leq k \leq n\}$  of  $d_1 \times d_2$  random matrices. Suppose that

$$\mathbb{E} \mathbf{S}_k = \mathbf{0} \quad \text{and} \quad \|\mathbf{S}_k\| \leq B \quad \text{for each index } k.$$

Introduce the sum of the random matrices:

$$\mathbf{Z} = \sum_{k=1}^n \mathbf{S}_k.$$

Define the matrix variance proxy:

$$\begin{aligned}\sigma^2 &= \max\{\|\mathbb{E} \mathbf{Z} \mathbf{Z}^*\|, \|\mathbb{E} \mathbf{Z}^* \mathbf{Z}\|\} \\ &= \max\left\{\left\|\sum_k \mathbb{E} \mathbf{s}_k \mathbf{s}_k^*\right\|, \left\|\sum_k \mathbb{E} \mathbf{s}_k^* \mathbf{s}_k\right\|\right\}.\end{aligned}$$

Then, for all  $t \geq 0$ ,

$$\mathbb{P}\{\|\mathbf{Z}\| \geq t\} \leq (d_1 + d_2) \cdot \exp\left(\frac{-t^2/2}{\sigma^2 + Bt/3}\right).$$

Furthermore,

$$\mathbb{E} \|\mathbf{Z}\| \leq \sqrt{2\sigma^2 \log(d_1 + d_2)} + \frac{1}{3}B \log(d_1 + d_2).$$

**Exercise 1.3** Establish Corollary 1.14 by applying Theorem 1.11 to the self-adjoint dilation  $\mathcal{H}(\mathbf{Z})$ , perhaps with larger constants. **Hint:** To obtain the sharp constants presented here, you need to use the maximum eigenvalue bound that appears inside the proof of Theorem 1.11.

## Notes

The modern theory of matrix concentration begins with the matrix Laplace transform technique (Proposition 1.1) developed by Ahlswede & Winter [AW02] and refined by Oliveira [Oli10]. The author of these notes recognized [Tro11a; Tro12; Tro15] that Lieb's theorem allows us to develop a perfect analogy (Theorem 1.7) with the scalar concentration theory. This idea has had a profound impact on computational mathematics over the last decade. These lectures explore some of the most striking outcomes.

Matrix concentration inequalities have a long history. Early work in operator theory and Banach space geometry includes [Buc01; Lus86; LP91; PX97; Rud99; Tom74]. The monograph [Tro15] provides a more comprehensive account.





## 2. Matrix Approximation by Sampling

“Corncobs” Wikimedia Commons

Most of the text in this lecture is copied from my monograph [Tro15, Chap. 6].

In applied mathematics, we often need to approximate a complicated target object by a more structured object. In some situations, we can solve this problem using a beautiful probabilistic approach called *empirical approximation*. The basic idea is to construct a “simple” random object whose expectation equals the target. We obtain the approximation by averaging several independent copies of the simple random object. As the number of terms in this average increases, the approximation becomes more complex, but it represents the target more faithfully. We must quantify this tradeoff.

In particular, we often encounter problems where we need to approximate a matrix by a more structured matrix. For example, we may wish to find a sparse matrix that is close to a given matrix, or we may need to construct a low-rank matrix that is close to a given matrix. Empirical approximation provides one mechanism for obtaining these approximations. The matrix Bernstein inequality offers a natural tool for assessing the quality of the randomized approximation.

This lecture develops a general framework for empirical approximation of symmetric matrices along with an application in machine learning. The monograph [Tro15, Chap. 6] includes the extension to rectangular matrices and several other basic applications.

### 2.1 Matrix sampling estimators

Let  $A$  be a self-adjoint target matrix that we hope to approximate by a more structured matrix. To that end, suppose we can represent the target as a sum of “simple” matrices:

$$A = \sum_{i=1}^N A_i. \quad (2.1)$$

The idea is to identify summands  $\mathbf{A}_i$  with desirable properties (such as sparsity or low rank) that we want our approximation to inherit.

Along with the decomposition (2.1), we need to construct a set of sampling probabilities:

$$\sum_{i=1}^N p_i = 1 \quad \text{and} \quad p_i > 0 \quad \text{for } i = 1, \dots, N. \quad (2.2)$$

We want to ascribe larger probabilities to “more important” summands. Quantifying what “important” means is the most difficult aspect of randomized matrix approximation. Choosing the right sampling distribution for a specific problem requires insight and ingenuity. Nevertheless, we will see that the matrix Bernstein inequality gives a strong hint about which distributions lead to the most accurate approximations.

Given the data (2.1) and (2.2), we may construct a “simple” random matrix  $\mathbf{R}$  by sampling:

$$\mathbf{R} = p_i^{-1} \mathbf{A}_i \quad \text{with probability } p_i. \quad (2.3)$$

This construction ensures that  $\mathbf{R}$  is an unbiased estimator of the target:  $\mathbb{E} \mathbf{R} = \mathbf{A}$ . Even so, the random matrix  $\mathbf{R}$  offers a poor approximation of the target  $\mathbf{A}$  because it has a lot more structure. To improve the quality of the approximation, we average  $n$  independent copies of the random matrix  $\mathbf{R}$ . We obtain an estimator of the form

$$\bar{\mathbf{R}}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{R}_k \quad \text{where each } \mathbf{R}_k \text{ is an independent copy of } \mathbf{R}.$$

By linearity of expectation, this estimator is also unbiased:  $\mathbb{E} \bar{\mathbf{R}}_n = \mathbf{A}$ . The approximation  $\bar{\mathbf{R}}_n$  remains structured when the number  $n$  of terms in the approximation is small as compared with the number  $N$  of terms in the decomposition (2.1).

Our goal is to quantify the approximation error as a function of the complexity  $n$  of the approximation:

$$\mathbb{E} \|\bar{\mathbf{R}}_n - \mathbf{A}\| \leq \text{err}(n).$$

As a reminder,  $\|\cdot\|$  denotes the spectral norm; i.e., the  $\ell_2$  operator norm. There is a tension between the total number  $n$  of terms in the approximation and the error  $\text{err}(n)$  the approximation incurs. In applications, it is essential to achieve the right balance.

### 2.1.1 An error estimate

We can obtain an error estimate for the approximation scheme described in Section 2.1 as an immediate corollary of the matrix Bernstein inequality.

**Theorem 2.1 (Matrix approximation by random sampling).** Let  $\mathbf{A} \in \mathbb{H}_d$  be a fixed matrix. Construct a random matrix  $\mathbf{R} \in \mathbb{H}_d$  that satisfies

$$\mathbb{E} \mathbf{R} = \mathbf{A} \quad \text{and} \quad \|\mathbf{R}\| \leq B.$$

Compute the per-sample second moment:

$$m_2(\mathbf{R}) = \|\mathbb{E} \mathbf{R}^2\|. \quad (2.4)$$



Form the matrix sampling estimator

$$\bar{\mathbf{R}}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{R}_k \quad \text{where each } \mathbf{R}_k \text{ is an independent copy of } \mathbf{R}.$$

Then the estimator satisfies, for all  $t \geq 0$ ,

$$\mathbb{P} \{ \|\bar{\mathbf{R}}_n - \mathbf{A}\| \geq t \} \leq 2d \exp \left( \frac{-nt^2/2}{m_2(\mathbf{R}) + 2Bt/3} \right). \quad (2.5)$$

Furthermore,

$$\mathbb{E} \|\bar{\mathbf{R}}_n - \mathbf{A}\| \leq \sqrt{\frac{2m_2(\mathbf{R}) \log(2d)}{n}} + \frac{2B \log(2d)}{3n}. \quad (2.6)$$

*Proof.* Since  $\mathbf{R}$  is an unbiased estimator of the target matrix  $\mathbf{A}$ , we can write

$$\mathbf{Y} = \bar{\mathbf{R}}_n - \mathbf{A} = \frac{1}{n} \sum_{k=1}^n (\mathbf{R}_k - \mathbb{E} \mathbf{R}) = \sum_{k=1}^n \mathbf{X}_k.$$

We have defined the summands  $\mathbf{X}_k = n^{-1}(\mathbf{R}_k - \mathbb{E} \mathbf{R})$ . These random matrices form an independent and identically distributed family, and each  $\mathbf{X}_k$  has mean zero.

Now, each of the summands is subject to an upper bound:

$$\|\mathbf{X}_k\| \leq \frac{1}{n} (\|\mathbf{R}_k\| + \|\mathbb{E} \mathbf{R}\|) \leq \frac{1}{n} (\|\mathbf{R}_k\| + \mathbb{E} \|\mathbf{R}\|) \leq \frac{2B}{n}.$$

The first relation is the triangle inequality; the second is Jensen's inequality. The last estimate follows from our assumption that  $\|\mathbf{R}\| \leq B$ .

To control the matrix variance, first note that

$$\left\| \sum_{k=1}^n \mathbb{E} \mathbf{X}_k^2 \right\| = n \cdot \|\mathbb{E} \mathbf{X}_1^2\|.$$

The identity holds because the summands  $\mathbf{X}_k$  are identically distributed. We may calculate that

$$\mathbf{0} \preceq \mathbf{X}_1^2 = n^{-2} \mathbb{E}(\mathbf{R} - \mathbb{E} \mathbf{R})^2 = n^{-2} [\mathbb{E} \mathbf{R}^2 - (\mathbb{E} \mathbf{R})^2] \preceq n^{-2} \mathbb{E} \mathbf{R}^2.$$

The first relation holds because the expectation of the random psd matrix  $\mathbf{X}_1^2$  is psd. The first identity follows from the definition of  $\mathbf{X}_1$  and the fact that  $\mathbf{R}_1$  has the same distribution as  $\mathbf{R}$ . The second identity is a direct calculation. The last relation holds because  $(\mathbb{E} \mathbf{R})^2$  is psd. In summary,

$$\left\| \sum_{k=1}^n \mathbb{E} \mathbf{X}_k^2 \right\| \leq \frac{1}{n} \|\mathbb{E} \mathbf{R}^2\| = \frac{m_2(\mathbf{R})}{n}.$$

The last line follows from the definition (2.4) of  $m_2(\mathbf{R})$ .

We are prepared to apply the matrix Bernstein inequality to the random matrix  $\mathbf{Y}$ . This act delivers the stated results.  $\blacksquare$

### 2.1.2 Discussion

One of the most common applications of the matrix Bernstein inequality is to analyze empirical matrix approximations. As a consequence, Corollary 2.1 is one of the most useful forms of the matrix Bernstein inequality. Let us discuss some of the important aspects of this result.

#### Understanding the bound on the approximation error

First, let us examine how many samples  $n$  suffice to bring the approximation error bound in Corollary 2.1 below a specified positive tolerance  $\varepsilon$ . Examining inequality (2.6), we find that

$$n \geq \frac{2m_2(\mathbf{R}) \log(2d)}{\varepsilon^2} + \frac{2B \log(2d)}{3\varepsilon} \quad \text{implies} \quad \mathbb{E} \|\bar{\mathbf{R}}_n - \mathbf{A}\| \leq 2\varepsilon. \quad (2.7)$$

Roughly, the number  $n$  of samples should be on the scale of the maximum of the per-sample second moment  $m_2(\mathbf{R})$  and the uniform upper bound  $B$ .

The bound (2.7) also reveals an unfortunate aspect of empirical matrix approximation. To make the tolerance  $\varepsilon$  small, the number  $n$  of samples must increase in proportion to  $\varepsilon^{-2}$ . In other words, it takes many samples to achieve a highly accurate approximation. We cannot avoid this phenomenon if we construct an approximation using an empirical average, because it is ultimately a consequence of the central limit theorem.

On a more positive note, it is quite valuable that the error bound (2.5) involves the spectral norm. This type of estimate simultaneously controls the error in every linear function of the approximation:

$$\|\bar{\mathbf{R}}_n - \mathbf{A}\| \leq \varepsilon \quad \text{implies} \quad |\text{tr}(\bar{\mathbf{R}}_n \mathbf{C}) - \text{tr}(\mathbf{A} \mathbf{C})| \leq \varepsilon \quad \text{for } \|\mathbf{C}\|_1 \leq 1.$$

We have written  $\|\cdot\|_1$  for the Schatten 1-norm. These bounds also control the error in each eigenvalue  $\lambda_j(\bar{\mathbf{R}}_n)$  of the approximation:

$$\|\bar{\mathbf{R}}_n - \mathbf{A}\| \leq \varepsilon \quad \text{implies} \quad |\lambda_j(\bar{\mathbf{R}}_n) - \lambda_j(\mathbf{A})| \leq \varepsilon.$$

When there is a gap between two eigenvalues of  $\mathbf{A}$ , we can also obtain bounds for the discrepancy between the associated eigenvectors of  $\bar{\mathbf{R}}_n$  and  $\mathbf{A}$  using perturbation theory [Bha97, Chap. VII].

#### Constructing empirical estimates

To obtain an accurate structured approximation, we need to select the right set of simple constituent matrices, as well as the right choice of sampling probabilities. In practice, these choices demand considerable creativity.

Fortunately, the matrix sampling result, Theorem 2.1, offers us some guidance because it identifies two summary parameters that control the quality of an empirical approximation. Indeed, we want to select the random matrix  $\mathbf{R}$  to ensure that the upper bound  $B$  and the per-sample second moment  $m_2(\mathbf{R})$  are both as small as possible. Later, we will see that this insight gives us a mechanism for determining the right sampling probabilities for certain problems.



This observation also hints at the possibility of achieving a bias–variance tradeoff when approximating  $\mathbf{A}$ . Indeed, we might drop all of the “unimportant” terms in the representation (2.1), i.e., those whose sampling probabilities are small. Then we construct a random approximation  $\mathbf{R}$  only for the “important” terms that remain. Properly executed, this process may decrease both the per-sample second moment  $m_2(\mathbf{R})$  and the upper bound  $B$ . The idea is analogous with shrinkage in statistical estimation.

### A general sampling model

Corollary 2.1 extends beyond the sampling model based on the finite expansion (2.1). Indeed, we can consider a general decomposition of the self-adjoint target matrix  $\mathbf{A}$ :

$$\mathbf{A} = \int_{\Omega} \mathbf{A}(\omega) d\mu(\omega), \quad (2.8)$$

where  $\mu$  is a probability measure on a sample space  $\Omega$ . As before, the idea is to represent the target matrix  $\mathbf{A}$  as an average of “simple” matrices  $\mathbf{A}(\omega)$ . The main difference is that the family of simple matrices may now be infinite. In this setting, we construct the random approximation  $\mathbf{R}$  so that

$$\mathbb{P}\{\mathbf{R} \in \mathbf{E}\} = \mu\{\omega : \mathbf{A}(\omega) \in \mathbf{E}\} \quad \text{for each Borel subset } \mathbf{E} \subseteq \mathbb{H}_d$$

In particular, it follows that

$$\mathbb{E} \mathbf{R} = \mathbf{A} \quad \text{and} \quad \|\mathbf{R}\| \leq \sup_{\omega \in \Omega} \|\mathbf{A}(\omega)\|.$$

In this lecture, we will see how this abstraction allows us to approximate kernel matrices for machine learning applications.

### Suboptimality of sampling estimators

Another fundamental point about sampling estimators is that they are often suboptimal. In other words, the matrix sampling estimator may incur an error substantially worse than the error in the best structured approximation of the target matrix.

To see why, let us consider a simple form of low-rank approximation by random sampling. The method here does not have practical value, but it highlights the reason that sampling estimators usually do not achieve ideal results. Suppose that  $\mathbf{A}$  is a trace-one psd matrix with the eigenvalue decomposition

$$\mathbf{A} = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^* \quad \text{where} \quad \sum_{i=1}^d \lambda_i = 1 \quad \text{and} \quad \lambda_i \geq 0.$$

Given the eigenvalue decomposition, we can construct a random rank-one approximation  $\mathbf{R}$  of the form

$$\mathbf{R} = \mathbf{u}_i \mathbf{u}_i^* \quad \text{with probability} \quad \lambda_i.$$

Per Corollary 2.1, the error in the associated sampling estimator  $\bar{\mathbf{R}}_n$  is a rank- $n$  matrix that satisfies

$$\|\bar{\mathbf{R}}_n - \mathbf{A}\| \leq \sqrt{\frac{2 \log(2d)}{n}} + \frac{2 \log(2d)}{n}$$

On the other hand, a best rank- $n$  approximation of  $\mathbf{A}$  takes the form  $\mathbf{A}_n = \sum_{j=1}^n \lambda_j \mathbf{u}_j \mathbf{u}_j^*$ , and it incurs error

$$\|\mathbf{A}_n - \mathbf{A}\| = \lambda_{n+1} \leq \frac{1}{n+1}.$$

The second relation is Markov's inequality, which provides an accurate estimate only when the singular values  $\lambda_1, \dots, \lambda_{n+1}$  are comparable. Regardless, the sampling estimator always incurs a somewhat larger error, which only converges as  $n^{-1/2}$ . Furthermore, there are many matrices whose singular values decay quickly, so that  $\lambda_{n+1} \ll (n+1)^{-1}$ . In the latter situation, the error in the sampling estimator is potentially much worse than the optimal error.

## 2.2 Application: Random features

As a first application of empirical matrix approximation, let us discuss an idea from machine learning called *random features*. The approach is based on the continuous sampling model (2.8), but it depends on the same principles as the discrete approximations that we introduced in Section 2.1.

Random feature maps were proposed by Ali Rahimi and Ben Recht [RR07], and they have turned out to be useful in practice. The analysis in this section is due to David Lopez-Paz et al. [Lop+14].

### 2.2.1 Kernel matrices

Let  $\mathcal{X}$  be a set. We think about the elements of the set  $\mathcal{X}$  as (potential) observations that we would like to use to perform learning and inference tasks. Let us introduce a bounded measure  $K$  of similarity between pairs of points in the set:

$$K : \mathcal{X} \times \mathcal{X} \rightarrow [-1, +1].$$

The similarity measure  $K$  is often called a *kernel*. We assume that the kernel returns the value +1 when its arguments are identical, and it returns smaller values when its arguments are dissimilar. We also assume that the kernel is symmetric; that is,  $K(\mathbf{x}, \mathbf{y}) = K(\mathbf{y}, \mathbf{x})$  for all arguments  $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ .

A simple example of a kernel is the angular similarity between a pair of points in a Euclidean space:

$$K(\mathbf{x}, \mathbf{y}) = \frac{2}{\pi} \arcsin \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|} = 1 - \frac{2}{\pi} \cdot \angle(\mathbf{x}, \mathbf{y}) \quad \text{for } \mathbf{x}, \mathbf{y} \in \mathbb{R}^d. \quad (2.9)$$

We write  $\angle(\cdot, \cdot)$  for the planar angle between two vectors, measured in radians. As usual, we instate the convention that  $0/0 = 0$ . See Figure 2.1 for an illustration.

Suppose that  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathcal{X}$  are observations. The  $N \times N$  kernel matrix  $\mathbf{G} = [g_{ij}]$  tabulates the values of the kernel function for each pair of data points:

$$g_{ij} = K(\mathbf{x}_i, \mathbf{x}_j) \quad \text{for } i, j = 1, \dots, N.$$

We say that the kernel  $K$  is *positive definite* if the kernel matrix  $\mathbf{G}$  is positive semidefinite for any choice of observations  $\{\mathbf{x}_i\} \subset \mathcal{X}$ . We will be concerned only with positive-definite kernels in this discussion. It may be helpful to think about the kernel matrix  $\mathbf{G}$  as a generalization of the Gram matrix of a family of points in a Euclidean space.

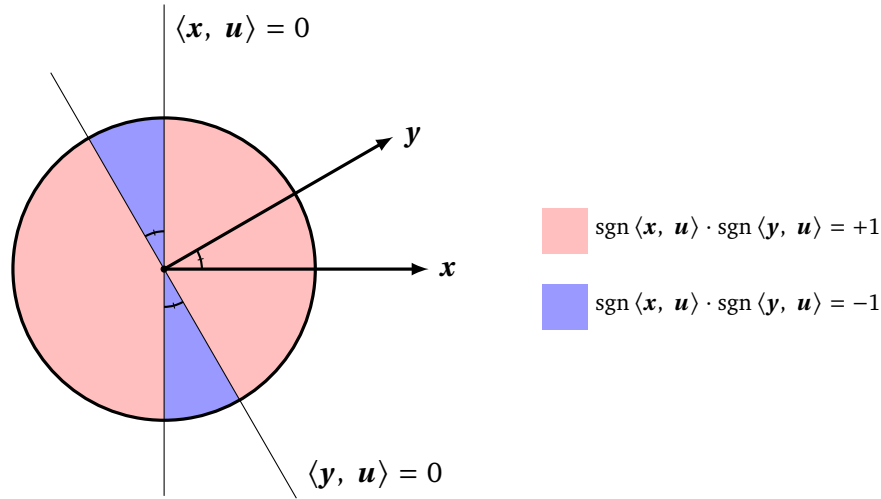


Figure 2.1: **The angular similarity between two vectors.** Let  $\mathbf{x}$  and  $\mathbf{y}$  be nonzero vectors in  $\mathbb{R}^2$  with angle  $\angle(\mathbf{x}, \mathbf{y})$ . The light red region contains the directions  $\mathbf{u}$  where the product  $\text{sgn } \langle \mathbf{x}, \mathbf{u} \rangle \cdot \text{sgn } \langle \mathbf{y}, \mathbf{u} \rangle$  equals  $+1$ , and the dark blue region contains the directions  $\mathbf{u}$  where the same product equals  $-1$ . The blue region subtends a total angle of  $2\angle(\mathbf{x}, \mathbf{y})$ , and the red region subtends a total angle of  $2\pi - 2\angle(\mathbf{x}, \mathbf{y})$ .

In the Euclidean setting, there are many statistical learning methods that only require the inner product between each pair of observations. These algorithms can be extended to the kernel setting by replacing each inner product with a kernel evaluation. As a consequence, kernel matrices can be used for classification, regression, and feature selection. In these applications, kernels are advantageous because they work outside the Euclidean domain, and they allow task-specific measures of similarity. This idea, sometimes called the *kernel trick*, is a major insight with wide applications [SS01].

A significant challenge for algorithms based on kernels is that the kernel matrix is big. Indeed,  $\mathbf{G}$  contains  $\Theta(N^2)$  entries, where  $N$  is the number of data points. Furthermore, the cost of constructing the kernel matrix is often  $\Theta(dN^2)$  where  $d$  is the number of parameters required to specify a point in the universe  $\mathcal{X}$ .

Nevertheless, there is an opportunity. Large data sets tend to be redundant, so the kernel matrix also tends to be redundant. This event manifests in the kernel matrix being well-approximated by a low-rank matrix. As a consequence, we may try to replace the kernel matrix by a low-rank proxy. For some similarity measures, we can accomplish this task using empirical approximation.

### 2.2.2 Random features and low-rank approximation of the kernel matrix

In certain cases, a positive-definite kernel can be written as an expectation (2.8), and we can take advantage of this representation to construct an empirical approximation of the kernel matrix. Let us begin with the general construction, and then we will present a few examples in Section 2.2.3.

Let  $\mathcal{W}$  be a sample space equipped with a sigma-algebra and a probability measure

$\mu$ . Introduce a bounded *feature map*:

$$\psi : \mathcal{X} \times \mathcal{W} \rightarrow [-b, +b] \quad \text{where } b > 0.$$

Consider a random variable  $\mathbf{w}$  taking values in  $\mathcal{W}$  and distributed according to the measure  $\mu$ . We assume that this random variable satisfies the *reproducing property*

$$K(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{\mathbf{w}} [\psi(\mathbf{x}; \mathbf{w}) \cdot \psi(\mathbf{y}; \mathbf{w})] \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathcal{X}. \quad (2.10)$$

The pair  $(\psi, \mathbf{w})$  is called a *random feature map* for the kernel  $K$ . As we will see, this hypothesis will lead to an instance of the expectation model (2.8) for the kernel matrix of an arbitrary dataset.

We want to approximate the kernel matrix associated with a set  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset \mathcal{X}$  of observations. To do so, we draw a random vector  $\mathbf{w} \in \mathcal{W}$  distributed according to  $\mu$ . Form a random vector  $\mathbf{z} \in \mathbb{R}^N$  by applying the feature map to each data point with the *same* choice of the random vector  $\mathbf{w}$ . That is,

$$\mathbf{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_N \end{bmatrix} = \begin{bmatrix} \psi(\mathbf{x}_1; \mathbf{w}) \\ \vdots \\ \psi(\mathbf{x}_N; \mathbf{w}) \end{bmatrix}.$$

The vector  $\mathbf{z} \in \mathbb{R}^N$  is sometimes called a *random feature*; it should be regarded as a summary of the entire dataset. By the reproducing property (2.10) for the random feature map, for each pair  $(i, j)$  of indices,

$$g_{ij} = K(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{E}_{\mathbf{w}} [\psi(\mathbf{x}_i; \mathbf{w}) \cdot \psi(\mathbf{x}_j; \mathbf{w})] = \mathbb{E}_{\mathbf{w}} [z_i \cdot z_j].$$

In other words, the feature map gives us an unbiased estimator for each entry of the kernel matrix.

We can write this relation in matrix form as

$$\mathbf{G} = \mathbb{E}[\mathbf{z}\mathbf{z}^*].$$

The random matrix  $\mathbf{R} = \mathbf{z}\mathbf{z}^*$  is an unbiased rank-one estimator for the kernel matrix  $\mathbf{G}$ . This is an instantiation of the model (2.8)! Note that this representation demonstrates that random feature maps, as defined here, only exist for positive-definite kernels. (But we can construct random feature maps for some other kinds of kernels using related approaches.)

We can construct a better empirical approximation of the kernel matrix  $\mathbf{G}$  by averaging realizations of the simple estimator  $\mathbf{R}$ :

$$\bar{\mathbf{R}}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{R}_k \quad \text{where } \mathbf{R}_k \text{ is an independent copy of } \mathbf{R}. \quad (2.11)$$

In other words, we are using  $n$  independent random features  $\mathbf{z}_1, \dots, \mathbf{z}_n$  to approximate the kernel matrix.

The cost of computing a single random feature is typically  $\Theta(dN)$ , where  $d$  is the number of parameters required to specify a point in the universe  $\mathcal{X}$ . Therefore, the cost of computing  $n$  random features is  $\Theta(dnN)$ . When  $n \ll N$ , the cost of obtaining the random feature approximation  $\bar{\mathbf{R}}_n$  is substantially smaller than the cost of computing the full kernel matrix. The question is how many random features  $n$  we needed before our estimator is accurate.

### 2.2.3 Examples of random feature maps

Before we continue with the analysis, let us describe some random feature maps. This discussion is tangential to our theme of matrix concentration, but it is valuable to understand why random feature maps exist.

#### The angular similarity kernel

First, let us consider the angular similarity (2.9) defined on  $\mathbb{R}^d$ . We can construct a random feature map using a classic result from plane geometry. If we draw  $\mathbf{w}$  uniformly from the unit sphere  $\mathbb{S}^{d-1} \subset \mathbb{R}^d$ , then

$$\begin{aligned} K(\mathbf{x}; \mathbf{y}) &= 1 - \frac{2}{\pi} \cdot \angle(\mathbf{x}, \mathbf{y}) \\ &= \mathbb{E}_{\mathbf{w}} \left[ \operatorname{sgn} \langle \mathbf{x}, \mathbf{w} \rangle \cdot \operatorname{sgn} \langle \mathbf{y}, \mathbf{w} \rangle \right] \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathcal{X}. \end{aligned} \quad (2.12)$$

The easy proof of this relation should be visible from the diagram in Figure 2.1. In light of the formula (2.12), we set  $\mathcal{W} = \mathbb{S}^{d-1}$  with the uniform measure, and we define the feature map

$$\psi(\mathbf{x}; \mathbf{w}) = \operatorname{sgn} \langle \mathbf{x}, \mathbf{w} \rangle.$$

The reproducing property (2.10) follows immediately from (2.12). Therefore, the pair  $(\psi, \mathbf{w})$  is a random feature map for the angular similarity kernel.

The paper [KK12] explains how to compute random features for more general inner-product kernels using a classic theorem of Schönberg.

#### Translation-invariant kernels

Next, let us describe an important class of kernels that can be expressed using random feature maps. A kernel on  $\mathbb{R}^d$  is *translation invariant* if there is a function  $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$  for which

$$K(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x} - \mathbf{y}) \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

Bôchner's theorem, a classical result from harmonic analysis, gives a representation for each continuous, positive-definite, translation-invariant kernel:

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= \varphi(\mathbf{x} - \mathbf{y}) \\ &= c \int_{\mathbb{R}^d} e^{i\langle \mathbf{x}, \mathbf{w} \rangle} \cdot e^{-i\langle \mathbf{y}, \mathbf{w} \rangle} d\mu(\mathbf{w}) \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^d. \end{aligned} \quad (2.13)$$

In this expression,  $c$  is a positive scale factor  $c$ , and  $\mu$  is a probability measure on  $\mathbb{R}^d$ , and these objects depend only on the function  $\varphi$ . Conversely, for any probability measure  $\mu$ , the formula (2.13) induces a continuous, positive-definite, translation-invariant kernel.

Bôchner's theorem (2.13) allows us to construct a (complex-valued) random feature map for the kernel  $K$ :

$$\psi_{\mathbb{C}}(\mathbf{x}; \mathbf{w}) = \sqrt{c} e^{i\langle \mathbf{x}, \mathbf{w} \rangle} \quad \text{where } \mathbf{w} \text{ has distribution } \mu \text{ on } \mathbb{R}^d.$$

This map satisfies a complex variant of the reproducing property (2.10):

$$K(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{\mathbf{w}} \left[ \psi_{\mathbb{C}}(\mathbf{x}; \mathbf{w}) \cdot \psi_{\mathbb{C}}(\mathbf{y}; \mathbf{w})^* \right] \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^d,$$

where we have written  $*$  for complex conjugation.

With a little more work, we can construct a real-valued random feature map. Recall that the kernel  $K$  is symmetric, so the complex exponentials in (2.13) can be written in terms of cosines. This observation leads to the random feature map

$$\psi(\mathbf{x}; \mathbf{w}, U) = \sqrt{2c} \cos(\langle \mathbf{x}, \mathbf{w} \rangle + U) \quad \text{where } \mathbf{w} \sim \mu \text{ and } U \sim \text{UNIFORM}[0, 2\pi]. \quad (2.14)$$

To verify that  $(\psi, (\mathbf{w}, U))$  reproduces the kernel  $K$ , as required by (2.10), we just make a short calculation using the angle-sum formula for the cosine.

We conclude this section with the most important example of a random feature map from the class we have just described. Consider the Gaussian radial basis function kernel:

$$K(\mathbf{x}, \mathbf{y}) = e^{-\alpha \|\mathbf{x} - \mathbf{y}\|^2 / 2} \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

The positive parameter  $\alpha$  reflects how close two points must be before they are regarded as “similar.” For the Gaussian kernel, Bôchner’s Theorem (2.13) holds with the scaling factor  $c = 1$  and the probability measure  $\mu = \text{NORMAL}(\mathbf{0}, \alpha \mathbf{I}_d)$ . In summary, we define

$$\psi(\mathbf{x}; \mathbf{w}, U) = \sqrt{2} \cos(\langle \mathbf{x}, \mathbf{w} \rangle + U) \quad \text{where } \mathbf{w} \sim \text{NORMAL}(\mathbf{0}, \alpha \mathbf{I}_d) \text{ and } U \sim \text{UNIFORM}[0, 2\pi].$$

This random feature map reproduces the Gaussian radial basis function kernel.

#### 2.2.4 Error bound for the random feature approximation

We will demonstrate that the approximation  $\bar{\mathbf{R}}_n$  of the  $N \times N$  kernel matrix  $\mathbf{G}$  using  $n$  random features, constructed in (2.11), leads to an estimate of the form

$$\mathbb{E} \|\bar{\mathbf{R}}_n - \mathbf{G}\| \leq \sqrt{\frac{2bN \|\mathbf{G}\| \log(2N)}{n}} + \frac{2bN \log(2N)}{3n}. \quad (2.15)$$

In this expression,  $b$  is the uniform bound on the magnitude of the feature map  $\psi$ . The short proof of (2.15) appears in Section 2.2.5.

To clarify what this result means, we introduce the *intrinsic dimension* of the  $N \times N$  kernel matrix  $\mathbf{G}$ :

$$\text{intdim}(\mathbf{G}) = \frac{\text{tr } \mathbf{G}}{\|\mathbf{G}\|} = \frac{N}{\|\mathbf{G}\|}.$$

Note that  $\text{tr } \mathbf{G} = N$  because of the requirement that  $K(\mathbf{x}, \mathbf{x}) = +1$  for all  $\mathbf{x} \in \mathcal{X}$ . The intrinsic dimension  $\text{intdim}(\mathbf{G})$  is a continuous measure of the number of energetic dimensions, and it is always bounded above by the algebraic rank of  $\mathbf{G}$ .

Now, assume that the number  $n$  of random features satisfies the bound

$$n \geq 2b\epsilon^{-2} \cdot \text{intdim}(\mathbf{G}) \cdot \log(2N),$$

In view of (2.15), the relative error in the empirical approximation of the kernel matrix satisfies

$$\frac{\mathbb{E} \|\bar{\mathbf{R}}_n - \mathbf{G}\|}{\|\mathbf{G}\|} \leq \epsilon + \epsilon^{-2}.$$

We learn that the randomized approximation of the kernel matrix  $\mathbf{G}$  using  $n$  random features can be accurate when  $n$  is proportional to the intrinsic dimension of  $\mathbf{G}$ , even if the intrinsic dimension is much smaller than the number of data points. That is,  $n \approx \text{intdim}(\mathbf{G}) \ll N$ .

### 2.2.5 Analysis of the random feature approximation

The analysis of random features is based on Corollary 2.1. To apply this result, we need the per-sample second-moment  $m_2(\mathbf{R})$  and the uniform upper bound  $B$ . Both are easy to come by.

First, observe that

$$\|\mathbf{R}\| = \|\mathbf{z}\mathbf{z}^*\| = \|\mathbf{z}\|^2 \leq bN$$

Recall that  $b$  is the uniform bound on the feature map  $\psi$ , and  $N$  is the number of components in the random feature vector  $\mathbf{z}$ .

Second, we calculate that

$$\mathbb{E} \mathbf{R}^2 = \mathbb{E} [\|\mathbf{z}\|^2 \mathbf{z}\mathbf{z}^*] \leq bN \cdot \mathbb{E}[\mathbf{z}\mathbf{z}^*] = bN \cdot \mathbf{G}.$$

Each random matrix  $\mathbf{z}\mathbf{z}^*$  is positive semidefinite, so we can introduce the upper bound  $\|\mathbf{z}\|^2 \leq bN$ . The last identity holds because  $\mathbf{R}$  is an unbiased estimator of the kernel matrix  $\mathbf{G}$ . It follows that

$$m_2(\mathbf{R}) = \|\mathbb{E} \mathbf{R}^2\| \leq bN \cdot \|\mathbf{G}\|.$$

This is our bound for the per-sample second moment.

Finally, we invoke Corollary 2.1 with parameters  $B = bN$  and  $m_2(\mathbf{R}) \leq bN\|\mathbf{G}\|$  to arrive at the estimate (2.15).







## 3. Quantum State Tomography

©CERN, CC BY-SA 3.0

This lecture was written primarily by Richard Kueng, on the basis of our joint work [Guh+18]. Any errors that appear are the fault of the lecturer.

A core problem in quantum information science is to estimate the state of a quantum system from measurements (of multiple realizations) of the system. This problem is called *quantum tomography*. In quantum computing, the state is represented by a finite-dimensional matrix, so we can formulate the tomography problem as a question about matrix estimation.

This lecture considers a special class of quantum tomography problems that admit a particularly simple analysis based on the matrix Bernstein inequality (Theorem 2.1). A remarkable feature of this application is that random matrices arise as a consequence of quantum mechanics!

### 3.1 Postulates of quantum mechanics

Quantum mechanics is a *probabilistic theory*, contra Einstein's firm belief that "God does not play dice." In this lecture, we will restrict ourselves to finite-dimensional quantum mechanics, where the principles are clearest. The extension to infinite dimensions is conceptually straightforward, and it resembles the transition from matrix analysis to functional analysis. To begin, we will develop the fundamental axioms of quantum mechanics as a noncommutative extension of discrete probability theory.

#### 3.1.1 Recapitulation: Discrete probability theory

Recall that  $\langle \cdot, \cdot \rangle$  is the standard inner product on  $\mathbb{R}^d$ . We denote the vector of ones by  $\mathbf{1} = (1, \dots, 1)^* \in \mathbb{R}^d$ ; this is the unit for the Hadamard product of vectors.

A discrete probability distribution on  $d$  points is fully characterized by a  $d$ -dimensional probability vector. A probability vector is just a nonnegative vector whose entries sum to one. The set  $\Delta_d$  of all  $d$ -dimensional probability vectors is called the *probability simplex*:

$$\Delta_d = \{\mathbf{p} \in \mathbb{R}^d : \mathbf{p} \geq \mathbf{0} \text{ and } \langle \mathbf{1}, \mathbf{p} \rangle = 1\}.$$

The probability simplex is a compact, convex set of vectors. The extreme points  $\delta_i$  of the probability simplex are the nonrandom probability distributions. The barycenter  $d^{-1}\mathbf{1}$  of the probability simplex is the uniform distribution.

Elementary events (singleton outcomes) are encoded by the  $d$  standard basis vectors  $\delta_1, \dots, \delta_d \in \mathbb{R}^d$ . Thus, the probability rule is given by the inner product:

$$\mathbb{P}\{i | \mathbf{p}\} = \langle \delta_i, \mathbf{p} \rangle = p_i \in [0, 1].$$

An event  $E$  is an element of the power set of  $\{1, \dots, d\}$ . We may represent the event  $E$  by the binary indicator vector  $\mathbf{1}_E \in \{0, 1\}^d$ . The probability rule remains the same:

$$\mathbb{P}\{E | \mathbf{p}\} = \langle \mathbf{1}_E, \mathbf{p} \rangle = \sum_{i \in E} p_i \in [0, 1].$$

This formalism extends to convex mixtures of events, which we call *generalized events*. The family of generalized events coincides with the standard cube:

$$\text{conv}\{0, 1\}^d = \{\mathbf{h} \in \mathbb{R}^d : \mathbf{0} \leq \mathbf{h} \leq \mathbf{1}\} = \mathbf{Q}_d.$$

The probability of a generalized event  $\mathbf{h} \in \mathbf{Q}_d$  is given by the inner product  $\langle \mathbf{h}, \mathbf{p} \rangle$ . We can therefore associate generalized events with the class of nonnegative random variables that are bounded by one, and the probability of a generalized event is the expectation of this random variable. In summary, generalized events are dual to probability distributions.

Next, we define the notion of a classical measurement.

**Definition 3.1 (Classical measurement).** A (classical) *measurement*  $\{\mathbf{h}_{\lambda_1}, \dots, \mathbf{h}_{\lambda_m}\} \subset \mathbf{Q}_d$  is a set of generalized events that forms a resolution of the vector of ones:

$$\mathbf{0} \leq \mathbf{h}_{\lambda_i} \leq \mathbf{1} \quad \text{and} \quad \sum_{i=1}^m \mathbf{h}_{\lambda_i} = \mathbf{1}.$$

A measurement should be viewed as a complete set of (generalized) events.

$$\sum_{i=1}^m \mathbb{P}\{\lambda_i | \mathbf{p}\} = \sum_{i=1}^m \langle \mathbf{h}_{\lambda_i}, \mathbf{p} \rangle = \langle \mathbf{1}, \mathbf{p} \rangle = 1.$$

In other words, it is certain that one of the outcomes  $\lambda_1, \dots, \lambda_m$  will occur.

**Example 3.1 (What is a classical measurement?).** Suppose that you and I agree on a bet that involves two random variables: a fair coin toss and the roll of a die. We first toss the coin and subsequently roll the die. The rules for victory depend on the outcome of the coin toss:

1. If the coin comes up heads, then I win if the die produces an odd number  $\{1, 3, 5\}$ . Otherwise, you win.

2. If the coin comes up tails, then I win if the die produces a number in the set  $\{1, 2, 3\}$ . Otherwise, you win.

A generalized classical event allows us to absorb the randomness in the coin flip into a generalized event that is associated only with the outcome of rolling the die:

$$\begin{aligned} \mathbf{h}_{\text{I win}} &= \frac{1}{2}(1, 0, 1, 0, 1, 0) + \frac{1}{2}(1, 1, 1, 0, 0, 0) = (1, 0.5, 1, 0, 0.5, 0); \\ \mathbf{h}_{\text{You win}} &= \frac{1}{2}(0, 1, 0, 1, 0, 1) + \frac{1}{2}(0, 0, 0, 1, 1, 1) = (0, 0.5, 0, 1, 0.5, 1). \end{aligned}$$

In other words, these two generalized events arise from marginalization over the first variable. The generalized events reflect our suspense about who will win the game once we roll the die.

The pair  $\{\mathbf{h}_{\text{I win}}, \mathbf{h}_{\text{You win}}\}$  constitutes a classical measurement system. In this case, performing the measurement amounts to completing the game (by rolling the die) and recording the outcome.

An alternative perspective, that is more quantum in spirit, realizes the probability rule as a tensor product and “sums out” one of the components. In the matrix setting, the analogous operation is called a partial trace. ■

Table 3.1 summarizes the basic concepts of classical discrete probability theory.

Concept	Representation	Formula
Probability density	Normalized, nonnegative $\mathbf{p} \in \mathbb{R}^d$	$\mathbf{p} \geq \mathbf{0}$ and $\langle \mathbf{1}, \mathbf{p} \rangle = 1$
Measurement	Resolution $\{\mathbf{h}_{\lambda_i}\}$ of the unit $\mathbf{1}$	$\mathbf{h}_{\lambda_i} \geq \mathbf{0}$ and $\sum_{i=1}^m \mathbf{h}_{\lambda_i} = \mathbf{1}$
Probability rule	Standard inner product	$\mathbb{P}\{\lambda_i   \mathbf{p}\} = \langle \mathbf{h}_{\lambda_i}, \mathbf{p} \rangle$

Table 3.1: **Axioms for classical probability theory.** The structure of discrete probability theory is captured by endowing  $\mathbb{R}^d$  with the partial order  $\geq$  and the identity element  $\mathbf{1}$ .

### 3.1.2 Noncommutative probability theory

The postulates of quantum mechanics arise naturally from a noncommutative extension of classical probability theory. We simply replace the triple  $(\mathbb{R}^d, \geq, \mathbf{1})$  by the triple  $(\mathbb{H}_d, \succcurlyeq, \mathbf{I})$ . Recall that  $\mathbb{H}_d = \mathbb{H}_d(\mathbb{C})$  is the space of self-adjoint  $d \times d$  complex matrices, endowed with the trace inner product  $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{tr}(\mathbf{XY})$ , the semidefinite order  $\succcurlyeq$ , and the identity matrix  $\mathbf{I}$ .

In quantum mechanics, the analog of a probability density vector is a (*probability*) *density matrix*.

**Definition 3.2 (Density matrix).** The state of a  $d$ -dimensional quantum mechanical system is fully described by a *density matrix*  $\rho \in \mathbb{H}_d$ , a self-adjoint matrix that satisfies

$$\rho \succcurlyeq \mathbf{0} \quad \text{and} \quad \langle \mathbf{I}, \rho \rangle = \text{tr}(\rho) = 1.$$

Density matrices are often called *states*.

Introduce the family  $\mathcal{S}(\mathbb{H}_d)$  of all  $d$ -dimensional density matrices:

$$\mathcal{S}(\mathbb{H}_d) = \{X \in \mathbb{H}_d : X \succcurlyeq \mathbf{0} \text{ and } \langle \mathbf{I}, X \rangle = 1\}$$

Like the probability simplex, the set  $\mathcal{S}(\mathbb{H}_d)$  of density matrices is compact and convex.

In parallel to a classical measurement system, we may now define a quantum measurement system.

**Definition 3.3 (Quantum measurement).** A (quantum) *measurement* is a collection  $\{H_{\lambda_i} : 1 \leq i \leq m\}$  of psd matrices that forms a resolution of the identity matrix:

$$\mathbf{0} \preccurlyeq H_{\lambda_i} \preccurlyeq \mathbf{I} \quad \text{and} \quad \sum_{i=1}^m H_{\lambda_i} = \mathbf{I}.$$

When a measurement  $\{H_{\lambda_i} : 1 \leq i \leq m\}$  is performed on a quantum mechanical system with density matrix  $\rho$ , two things happen.

1. **Born's rule:** We obtain a random measurement outcome  $\lambda_i$  that follows the probability distribution

$$\mathbb{P}\{\lambda_i | \rho\} = \langle H_{\lambda_i}, \rho \rangle = \text{tr}(H_{\lambda_i} \rho). \quad (3.1)$$

2. **Collapse of wavefunction:** The quantum system ceases to exist.

There is some philosophical debate about this model, but experimental evidence suggests that it serves well as an ideal representation of what happens in real quantum systems.

Table 3.2 summarizes the essential concepts of quantum probability theory. We remark that the transition from classical to quantum probability theory resembles the transition from linear to semidefinite programming.

Concept	Representation	Formula
Probability density	Normalized psd matrix $\rho \in \mathbb{H}_d$	$\rho \succcurlyeq \mathbf{0}$ and $\langle \mathbf{I}, \rho \rangle = 1$
Measurement	Resolution $\{H_{\lambda_i}\}$ of the identity $\mathbf{I}$	$H_{\lambda_i} \succcurlyeq \mathbf{0}$ and $\sum_{i=1}^m H_{\lambda_i} = \mathbf{I}$
Born's rule	Trace inner product	$\mathbb{P}\{\lambda_i   \rho\} = \langle H_{\lambda_i}, \rho \rangle$

Table 3.2: **Axioms for quantum mechanics.** The structure of quantum mechanics is captured by the real-linear space  $\mathbb{H}_d$  endowed with the psd order  $\succcurlyeq$  and the identity matrix  $\mathbf{I}$ .

### 3.1.3 Aside: Geometric intuition and the Bloch ball

Since the set  $\mathcal{S}(\mathbb{H}_d)$  of states is a convex body, we can distinguish points that capture information about its geometry.

- **Extreme points:** A density matrix  $\rho$  is an extreme point of  $\mathcal{S}(\mathbb{H}_d)$  if and only if  $\rho$  has rank one. Equivalently,  $\rho = uu^*$  where  $u \in \mathbb{C}^d$  is a unit vector. Extreme points are called *pure (quantum) states*, and they generate  $\mathcal{S}(\mathbb{H}_d)$  via convex mixtures:

$$\mathcal{S}(\mathbb{H}_d) = \text{conv}\{uu^* : u \in \mathbb{C}^d \text{ and } \|u\| = 1\}.$$

A pure state is the quantum analog of a classical nonrandom distribution.

- **Barycenter:** The barycenter of  $\mathcal{S}(\mathbb{H}_d)$  is the state  $\rho_0 = d^{-1}\mathbf{I}$ . It is called the *maximally mixed (quantum) state*. The maximally mixed state is the quantum analog of the classical uniform distribution.

For two-dimensional quantum states (called *qubits*), we can construct a beautiful geometric representation, called the *Bloch ball*. This representation helps us visualize the structure of the set of qubits, including the relationships between pure states and the maximally mixed state.

To construct the Bloch ball, we first define the *Pauli matrices*:

$$\sigma_0 = \mathbf{I}; \quad \sigma_1 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}; \quad \sigma_2 = \begin{bmatrix} 0 & -i \\ i & 0 \end{bmatrix}; \quad \sigma_3 = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}.$$

It is straightforward to check that these matrices form a basis of  $\mathbb{H}_2$ . Consider the parameterized family

$$\mathbf{M}(\mathbf{r}) = \sum_{i=0}^4 r_i \sigma_i \quad \text{where } \mathbf{r} = (r_0, r_1, r_2, r_3)^* \in \mathbb{R}^4.$$

We can easily characterize when  $\mathbf{M}(\mathbf{r})$  is a density matrix:

- $\mathbf{M}(\mathbf{r})$  has unit trace if and only if  $r_0 = \frac{1}{2}$ .
- $\mathbf{M}(\mathbf{r})$  is psd if and only if  $r_1^2 + r_2^2 + r_3^2 \leq r_0^2$ .

In other terms,

$$\mathcal{S}(\mathbb{H}_2) = \left\{ \frac{1}{2} \sigma_0 + \frac{1}{2} \sum_{i=1}^3 r_i \sigma_i : r_1^2 + r_2^2 + r_3^2 \leq 1 \right\}. \quad (3.2)$$

The set of qubits is parameterized by linear combinations of Pauli matrices whose expansion coefficients  $\mathbf{r}' = (r_1, r_2, r_3) \in \mathbb{R}^3$  are confined to the unit ball.

The formula (3.2) establishes a one-to-one correspondence between the density matrices of two-dimensional quantum systems and the Euclidean unit ball  $\mathbb{S}^2 \subset \mathbb{R}^3$ . This is called the *Bloch ball representation*, and it accurately reflects the geometry of  $\mathcal{S}(\mathbb{H}_2)$ . Indeed,

- The maximally mixed state  $\rho_0 = \frac{1}{2}\mathbf{I}$  is associated with the point  $\mathbf{r}' = \mathbf{0} \in \mathbb{R}^3$ , the center of the Bloch ball.
- A density matrix  $\rho$  is a pure state if and only if the associated vector  $\mathbf{r}' \in \mathbb{R}^3$  of expansion coefficients has unit norm. This observation establishes a one-to-one correspondence between pure states (the extreme points of  $\mathbb{H}_2$ ) and unit vectors (the extreme points of  $\mathbb{S}^2$ ).

Figure 3.1 contains an illustration of this correspondence.

**Example 3.2 (Stern–Gerlach experiment).** Depending on the measurement, a single density matrix can produce both a completely deterministic and a uniformly random outcome distribution. This observation is at the heart of the famous Stern–Gerlach experiment (1922), one of the first demonstrations of genuine quantum behavior.



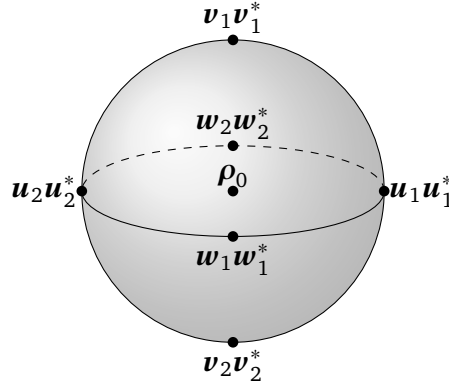


Figure 3.1: **Bloch ball representation of  $\mathcal{S}(\mathbb{H}^2)$ .** The maximally mixed state  $\rho_0 = \frac{1}{2}\mathbf{I}$  lies at the center of the Bloch ball. The surface of the ball is in one-to-one relation with the set of all pure quantum states. Also displayed: Three pairs of mutually orthogonal pure states that are evenly distributed across the boundary of  $\mathcal{S}(\mathbb{H}^2)$ .

Fix  $d = 2$ . Define four unit vectors:

$$\begin{aligned} \mathbf{u}_1 &= (1, 0)^* & \text{and} & & \mathbf{v}_1 &= \frac{1}{\sqrt{2}}(1, 1)^*; \\ \mathbf{u}_2 &= (0, 1)^* & \text{and} & & \mathbf{v}_2 &= \frac{1}{\sqrt{2}}(1, -1)^*. \end{aligned}$$

Then  $\{\mathbf{u}_1 \mathbf{u}_1^*, \mathbf{u}_2 \mathbf{u}_2^*\}$  and  $\{\mathbf{v}_1 \mathbf{v}_1^*, \mathbf{v}_2 \mathbf{v}_2^*\}$  describe two different quantum measurements. Applied to the pure two-dimensional state  $\rho = \mathbf{u}_1 \mathbf{u}_1^*$ , these measurements yield radically different outcome distributions:

$$\text{Measurement I: } \mathbb{P}\{1 | \rho\} = |\langle \mathbf{u}_1, \mathbf{u}_1 \rangle|^2 = 1 \quad \text{and} \quad \mathbb{P}\{2 | \rho\} = 0.$$

$$\text{Measurement II: } \mathbb{P}\{2 | \rho\} = |\langle \mathbf{v}_1, \mathbf{u}_1 \rangle|^2 = \frac{1}{2} \quad \text{and} \quad \mathbb{P}\{1 | \rho\} = \frac{1}{2}.$$

We refer to Figure 3.1 for a visualization of the underlying geometry. ■

### 3.2 Quantum state tomography

*Quantum state tomography* is the task of reconstructing the density matrix of a quantum system from measurement data. Quantum tomography is one of the oldest and most fundamental learning problems in quantum information science. Today, quantum tomography is a routine task that is essential for designing, testing, and tuning qubits in our quest to building scalable devices for quantum information processing.

Recall that a density matrix  $\rho \in \mathcal{S}(\mathbb{H}_d)$  contains a complete description of a  $d$ -dimensional quantum system. Knowledge of the density matrix, therefore, allows us to make predictions about future quantum measurements of an equivalent system. It also contains information about quantum-mechanical aspects of the system. For example, we can compute the *purity* of the system (i.e., the approximate rank of  $\rho$ ) and the *entanglement* among subsystems of a multipart system (i.e., how strongly the subsystems are correlated).

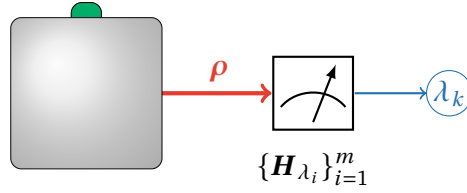


Figure 3.2: **Schematic of quantum state tomography.** A black box (*machine*) that is capable of producing a quantum system with density matrix  $\rho$  upon request (red arrow). A subsequent quantum measurement (*gauge*) yields a single outcome  $\lambda_k$  (blue arrow), but destroys the quantum system. The procedure must be repeated on fresh copies of the state in order to obtain additional information. Quantum state tomography is the task of reconstructing  $\rho$  from multiple observed outcomes.

In most settings, the density matrix  $\rho$  is not directly accessible. Instead, we obtain indirect information by performing a quantum measurement  $\{\mathbf{H}_{\lambda_i} : 1 \leq i \leq m\}$ . Born's rule (3.1) asserts that data about  $\rho$  is encoded in the probability distribution of outcomes (rather than the specific outcome  $\lambda_k$ ). Unfortunately, after a measurement is performed, the quantum system ceases to exist. To counter this challenge, we can prepare many copies of the same state, measure each one independently, and combine the information to estimate the distribution of outcomes accurately. Figure 3.2 contains a schematic.

Mathematically, this estimation problem combines interesting aspects of several scientific disciplines, most notably geometry and statistics.

### 3.2.1 Geometric aspects and measurement design

To build up some intuition, we first ignore the statistical aspects of quantum state tomography. Let  $\{\mathbf{H}_{\lambda_i}\}_{i=1}^m \subset \mathbb{H}_d$  be a fixed measurement. Suppose that we have the capacity to repeatedly perform this measurement on  $n$  realizations of an unknown quantum state  $\rho \in \mathcal{S}(\mathbb{H}_d)$ , where  $n \rightarrow \infty$ . In principle, this operation would allow us to determine the exact (classical) distribution of outcomes:

$$\mathbf{p} \in \Delta_m \quad \text{where} \quad p_k = \mathbb{P}\{\lambda_k \mid \rho\} = \langle \mathbf{H}_{\lambda_k}, \rho \rangle \quad \text{for } 1 \leq k \leq m. \quad (3.3)$$

Thus, Born's rule (3.1) describes a linear map (3.3) between the set  $\mathcal{S}(\mathbb{H}_d)$  of density matrices and the set  $\Delta_m \subset \mathbb{R}^m$  of classical probability distributions.

In this mathematical idealization, quantum state tomography becomes a linear inverse problem: Recover  $\rho \in \mathcal{S}(\mathbb{H}_d)$  from its linear image  $\mathbf{p} \in \Delta_m$ . This task is possible if and only if the linear measurement map (3.3) is injective. The following definition captures this idea.

**Definition 3.4 (Tomographic completeness).** A quantum measurement system  $\{\mathbf{H}_{\lambda_i} : 1 \leq i \leq m\} \subset \mathbb{H}_d$  is *tomographically complete* if and only if, for each pair  $\rho, \sigma \in \mathcal{S}(\mathbb{H}_d)$  of distinct states, there exists an index  $k \in \{1, \dots, m\}$  such that  $\langle \mathbf{H}_{\lambda_k}, \rho \rangle \neq \langle \mathbf{H}_{\lambda_k}, \sigma \rangle$ .



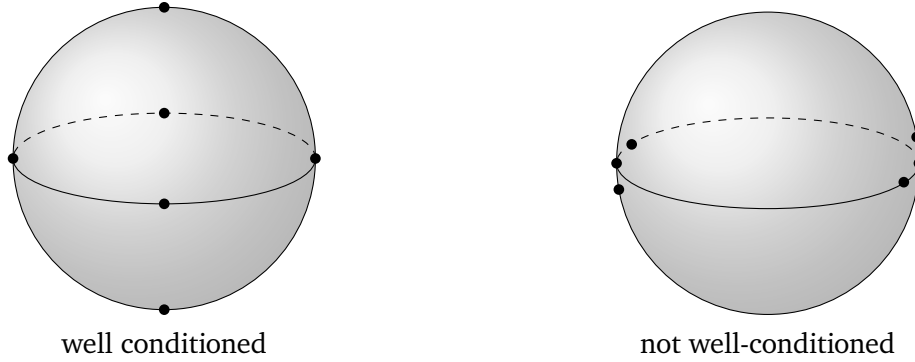


Figure 3.3: **Bloch sphere representation of two tomographically complete measurements in  $\mathbb{H}_2$ .** [left] The elements (*points*) of the measurement system are spread out evenly, which results in a well-conditioned measurement map. Indeed, it is a pair of mutually unbiased bases. [right] The elements (*points*) of the measurement system cluster at opposite extremes. This measurement is ill-equipped to accurately resolve points on the Bloch sphere in the vicinity of the north and south poles. As a consequence, the associated measurement is not well-conditioned.

A measurement cannot be tomographically complete unless it contains a sufficiently large number of outcomes. Indeed, the number  $m$  of potential measurement outcomes must obey

$$m \geq \dim \mathcal{S}(\mathbb{H}_d) + 1 = \dim \mathbb{H}_d = d^2.$$

This is just a basic fact about linear algebra.

Tomographic completeness is not the only property that we require of a measurement system. Indeed, injectivity only implies that the condition number<sup>1</sup> of the measurement map (3.3) is finite. If the condition number  $\kappa$  is large, we will suffer large errors when we try to solve the inverse problem (from a finite amount of data). In contrast, when the condition number  $\kappa \approx 1$ , the inverse problem can be solved in a stable fashion. Refer to Figure 3.3 for an illustration of two extreme cases, via the Bloch ball representation.

A linear map has the minimal condition number  $\kappa = 1$  if and only if the linear map is an isometry. Unfortunately, a quantum measurement map (3.3) can never be tomographically complete and isometric at the same time! The following definition describes the best-conditioned measurement maps that do exist.

**Definition 3.5 (Near-isotropic quantum measurement).** A quantum measurement system  $\{H_{\lambda_i} : 1 \leq i \leq m\} \subset \mathbb{H}_d$  is *near isotropic* when

1. Each element  $H_{\lambda_i} = (d/m) \mathbf{v}_i \mathbf{v}_i^*$  where  $\mathbf{v}_i \in \mathbb{C}^d$  is a unit vector;
2. The measurement has the reconstruction property

$$\frac{1}{m} \sum_{i=1}^m \langle \mathbf{v}_i \mathbf{v}_i^*, \mathbf{X} \rangle \mathbf{v}_i \mathbf{v}_i^* = \frac{1}{(d+1)d} (\mathbf{X} + (\text{tr } \mathbf{X}) \cdot \mathbf{I}) \quad \text{for all } \mathbf{X} \in \mathbb{H}_d. \quad (3.4)$$

<sup>1</sup>The condition number of a linear map is the ratio between largest and smallest singular value.

The linear map (3.3) associated with a near-isotropic quantum measurement has condition number  $\kappa$  that is bounded independent of the state dimension  $d$ . Moreover, the condition number  $\kappa \rightarrow 1$  as the state dimension  $d \rightarrow \infty$ .

In other research areas, you may encounter a system  $\{\mathbf{v}_1, \dots, \mathbf{v}_m\} \subset \mathbb{C}^d$  of unit vectors that satisfies (3.4). In approximation theory, these systems are interpreted as quadrature rules for polynomials on the complex unit sphere, and they are known as (complex projective) 2-designs. In frame theory, these systems are called tight fusion frames.

There are many interesting constructions of near-isotropic quantum measurements that arise from these connections.

**Example 3.3 (Near-isotropic quantum measurements).** The following quantum measurement systems are near-isotropic.

1. The *uniform measurement* is the infinite family  $\{d \mathbf{v} \mathbf{v}^*\}$  of all rescaled rank-one projectors  $d \mathbf{v} \mathbf{v}^*$ , endowed with the unique rotation-invariant probability measure  $d\mathbf{v}$  on the complex unit sphere.
2. The union of  $d + 1$  mutually unbiased bases<sup>2</sup> forms a set of  $m = (d + 1)d$  unit vectors that obey (3.4). Explicit constructions of these families are known when the dimension  $d$  is a prime power.
3. A set of  $m = d^2$  equiangular lines in  $\mathbb{C}^d$  also obeys (3.4). Zauner's conjecture states that such configurations should exist in any dimension  $d$ . A complete proof of this conjecture is an important open problem.

In theory, each of these measurement systems provides an effective way of acquiring information about a quantum state. But these measurements are not always realizable in practical settings. Constructing well-conditioned, implementable measurements is a major challenge in quantum engineering. ■

### 3.2.2 Statistical aspects and convergence

So far, we have considered a mathematically ideal version of quantum state tomography in which we measure an infinite number of realizations of the same state. In practice, the number  $n$  of samples is necessarily finite, so we cannot determine the probability distribution over the measurement outcomes exactly. As a consequence, we cannot expect to recover an unknown density matrix exactly.

Nevertheless, we can obtain an accurate approximation of the state with high probability, provided that we can perform a sufficient number of measurements. To assess the accuracy of an estimate  $\hat{\rho} \in \mathbb{H}_d$  of a state  $\rho \in \mathcal{S}(\mathbb{H}_d)$ , we typically use the Schatten 1-norm  $\|\hat{\rho} - \rho\|_1$ . The Schatten 1-norm, also known as the *trace norm*, is the quantum analog of the total-variation distance that arises in classical probability. Furthermore, this error measure has a natural operational interpretation in terms of quantum hypothesis testing.

**Definition 3.6 (Sample complexity).** Fix parameters  $\varepsilon, \delta \in (0, 1)$  and a rank  $1 \leq r \leq d$ . Let  $\rho \in \mathbb{H}_d$  be an unknown state with rank  $r$ . Perform the same quantum measurement

<sup>2</sup>Two orthonormal bases  $\{\mathbf{b}_1, \dots, \mathbf{b}_d\}$  and  $\{\mathbf{c}_1, \dots, \mathbf{c}_d\} \subset \mathbb{C}^d$  are mutually unbiased if  $|\langle \mathbf{b}_i, \mathbf{c}_j \rangle|^2 = d^{-1}$  for all  $1 \leq i, j \leq d$ . The standard basis and the discrete Fourier basis provide an instructive example.

on  $n$  realizations of the state, and construct a tomographic estimate  $\hat{\rho}_n \in \mathbb{H}_d$ . The *sample complexity* of this family of estimators is the minimum number  $n$  required to estimate the state with high probability:

$$\|\hat{\rho}_n - \rho\|_1 \leq \varepsilon \quad \text{with probability } 1 - \delta.$$

The sample complexity will typically depend on the rank  $r$ .

Methods from quantum information theory lead to a rigorous lower bound on the sample complexity of *any* tomographic estimation procedure [Haa+17].

**Theorem 3.7 (Haah et al. 2017 — informal).** Any tomographic estimator based on repeating the same measurement has sample complexity

$$n \gtrsim r^2 d \varepsilon^{-2} \log(1/\delta).$$

### 3.3 Quantum state tomography via matrix sampling

We are going to present and analyze a simple, yet powerful, estimation technique for quantum state tomography based on matrix sampling. To motivate the approach, we first consider a problem in classical probability.

#### 3.3.1 Estimating the bias of a coin

A classical analog of quantum state estimation is the problem of estimating the bias in a coin by flipping it repeatedly.

A coin is a two-dimensional classical random variable that is described by a single parameter, the *bias*  $p \in [0, 1]$ . The two outcomes follow the distribution  $\mathbb{P}\{\text{heads}\} = p$  and  $\mathbb{P}\{\text{tails}\} = 1 - p$ . How can we estimate the bias  $p$  by repeatedly tossing the coin? The simplest approximation procedure is based on a simple and intuitive decision rule. Toss the coin once and set

$$\hat{p} = \begin{cases} 1, & \text{if heads;} \\ 0, & \text{if tails.} \end{cases}$$

In general, this is a terrible estimator. But it does have the virtue of being unbiased:

$$\mathbb{E} \hat{p} = p \times 1 + (1 - p) \times 0 = p.$$

Instead, we toss the coin  $n$  times, form the estimators  $\hat{p}_1, \dots, \hat{p}_n$ , and construct the empirical average:

$$\bar{p}_n = n^{-1} \sum_{i=1}^n \hat{p}_i.$$

The empirical average will converge to the true bias of the coin.

To verify this claim and obtain a convergence rate, just apply Theorem 2.1 with  $d = 1$ . For each  $t \in [0, 1]$ ,

$$\mathbb{P}\{|\bar{p}_n - p| \geq t\} \leq 2 \exp\left(\frac{-nt^2/2}{p + 2t/3}\right) \leq 2 \exp\left(-\frac{3}{10}nt^2\right).$$

Therefore, for any parameters  $\varepsilon, \delta \in (0, 1)$ , if we perform

$$n \geq \frac{10}{3} \varepsilon^{-2} \log(1/(2\delta))$$

independent coin tosses, then the sample average satisfies the error bound  $|\bar{\rho}_n - \rho| < \varepsilon$  with probability at least  $1 - \delta$ .

### 3.3.2 The matrix sampling estimator

The simple coin tossing example can readily be generalized to quantum state tomography in  $d \geq 2$  dimensions. Construct a quantum system with unknown density matrix  $\rho \in \mathcal{S}(\mathbb{H}_d)$ , and perform the near-isotropic quantum measurement

$$\{H_{\lambda_i} = (d/m) \mathbf{v}_i \mathbf{v}_i^* : 1 \leq i \leq m\}.$$

When reading the measurement outcome, set

$$\mathbf{R} = \begin{cases} (d+1) \mathbf{v}_1 \mathbf{v}_1^* - \mathbf{I}, & \text{if we observe outcome } \lambda_1; \\ \vdots & \\ (d+1) \mathbf{v}_m \mathbf{v}_m^* - \mathbf{I}, & \text{if we observe outcome } \lambda_m. \end{cases} \quad (3.5)$$

Born's rule (3.1) and the geometric properties of near-isotropic measurements (3.4) ensure that the quantum estimator (3.5) is correct in expectation:

$$\begin{aligned} \mathbb{E}(\mathbf{R} + \mathbf{I}) &= \sum_{i=1}^m \mathbb{P}\{\lambda_i | \rho\} \cdot (d+1) \mathbf{v}_i \mathbf{v}_i^* \\ &= \frac{(d+1)d}{m} \sum_{i=1}^m \text{tr}(\mathbf{v}_i \mathbf{v}_i^* \rho) \mathbf{v}_i \mathbf{v}_i^* = \rho + \mathbf{I}. \end{aligned} \quad (3.6)$$

We also remark that  $\mathbf{R}$  has trace one, but it need not be psd.

We repeat this estimation procedure  $n$  times, for  $n$  copies of the quantum system, and we construct the sample average:

$$\bar{\mathbf{R}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{R}_i \quad \text{where } \mathbf{R}_i \text{ are iid copies of } \mathbf{R}.$$

The sample average has trace one, and it is an unbiased estimator of the state. On the other hand, it is not always psd. See Figure 3.4 for an illustration of the convergence of this sequence of estimates.

This construction is formally similar to the one in Lecture 2, but let us point out a major conceptual difference. Before, we designed an algorithm that makes random choices to construct random matrices that approximate a kernel matrix. In quantum state tomography, the estimator (3.5) produces independent random matrices because of the laws of quantum mechanics.

### 3.3.3 Sample complexity of the sample average

We quickly derive the convergence rate of the sample average using Theorem 2.1.

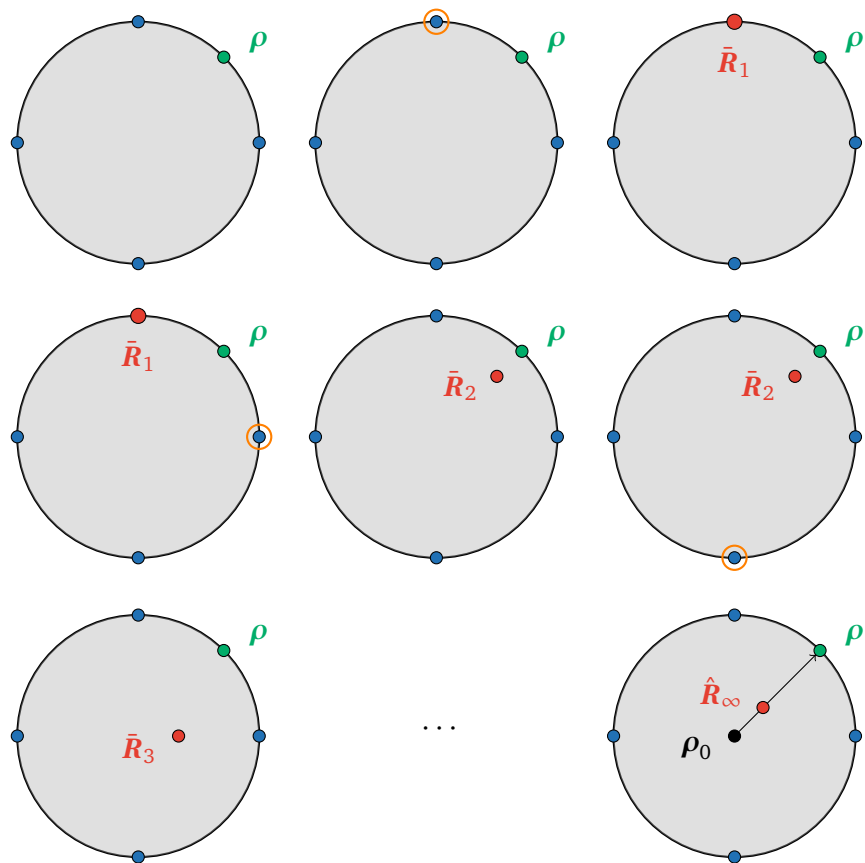


Figure 3.4: **Convergence of a naïve tomography estimator.** Perform a near-isotropic measurement (*blue points*) on an unknown density matrix  $\rho$  (*green*). Upon receiving a certain outcome (*orange circle*), we estimate  $\rho$  by the measurement element associated with this outcome (*red*). Repeat this procedure  $n$  times and construct the empirical average  $\bar{R}_n$ . Convex mixing properties imply that this estimator is pushed inside the set of quantum states (*convex combination*). But it maintains the correct direction in the limit of many repetitions. This illustration is exact if we restrict attention to the equatorial plane of the Bloch ball representation of  $\mathcal{S}(\mathbb{H}_2)$ .

**Proposition 3.8 (Quantum state tomography with sample averages).** Suppose that we have access to  $n$  (unentangled) realizations of a quantum system with density matrix  $\rho \in \mathcal{S}(\mathbb{H}_d)$ . Perform a near-isotropic quantum measurement on the  $i$ th realization, and construct the estimator  $\mathbf{R}_i \in \mathbb{H}_d$  based on the decision rule (3.5). Then the *matrix sample average*  $\bar{\mathbf{R}}_n = n^{-1} \sum_{i=1}^n \mathbf{R}_i$  obeys the error estimate

$$\mathbb{P} \{ \|\bar{\mathbf{R}}_n - \rho\| \geq t \} \leq 2d \exp \left( \frac{-3nt^2}{16d} \right).$$

This formula is valid for all  $t \in [0, 1]$ .

*Proof.* Let  $\mathbf{R}$  be the estimator (3.5). First, compute the upper bound

$$B = \sup \|\mathbf{R}\| = \max_{i=1, \dots, m} \|(d+1)\mathbf{v}_i\mathbf{v}_i^* - \mathbf{I}\| = d.$$

Another short calculation shows that the per-sample second moment satisfies

$$m_2(\mathbf{R}) = \|\mathbb{E} \mathbf{R}^2\| = \|(d-1)\rho + d\mathbf{I}\| \leq 2d.$$

We leave the details as an easy exercise.

Now, quantum measurements of unentangled quantum systems ensures are statistically independent. Therefore, the random matrices  $\mathbf{R}_i$  are independent copies of the random matrix  $\mathbf{R}$ . Theorem 2.1 implies that the matrix sample estimator  $\bar{\mathbf{R}}_n = n^{-1} \sum_{i=1}^n \mathbf{R}_i$  concentrates sharply around its expectation  $\rho$ . For  $t \geq 0$ ,

$$\mathbb{P} \{ \|\bar{\mathbf{R}}_n - \rho\| \geq t \} \leq 2d \exp \left( \frac{-nt^2/2}{m_2(\mathbf{R}) + 2Bt/3} \right) \leq 2d \exp \left( \frac{-3nt^2}{16d} \right).$$

This is what we needed to show. ■

### 3.3.4 Projection onto the set of quantum states

Proposition 3.8 equips the matrix sample average estimator  $\bar{\mathbf{R}}_n$  with a rigorous guarantee that it converges to the unknown density matrix  $\rho \in \mathcal{S}(\mathbb{H}_d)$ . For  $\tau \in (0, 1)$ , a total of  $n \gtrsim d \log(d)/\tau^2$  measurement repetitions are sufficient to ensure that  $\|\bar{\mathbf{R}}_n - \rho\| \leq \tau$  with high probability. Although powerful, this statement has two drawbacks:

1. The matrix sample estimator  $\bar{\mathbf{R}}_n$  is typically not psd. We therefore estimate the state  $\rho$  by something that is not itself a state.
2. Accuracy is reported in operator norm distance, rather than trace-norm distance.

Surprisingly, both drawbacks can be overcome by a single refinement. Just replace the sample matrix estimator by the closest density matrix, computed with respect to the Frobenius norm:

$$\hat{\rho}_n = \arg \min_{\sigma \in \mathcal{S}(\mathbb{H}_d)} \|\sigma - \bar{\mathbf{R}}_n\|_F. \quad (3.7)$$

We call this estimator the *projected matrix sample average*.

Intuitively, the projection onto the density matrices should decrease the distance between the estimator and target state. The following technical result makes this claim precise.

**Lemma 3.9** Fix a rank- $r$  density matrix  $\rho \in \mathcal{S}(\mathbb{H}_d)$  and a matrix  $M \in \mathbb{H}_d$  with trace one. Then the closest density matrix  $\sigma$  to  $M$  necessarily obeys

$$\|\rho - \sigma\|_1 \leq 4r \|\rho - M\|.$$

*Proof sketch.* The difference  $X = \rho - \sigma$  is a traceless self-adjoint matrix. Moreover, the positive part of  $X$  has rank no greater than  $r$  because both  $\rho$  and  $\sigma$  are psd and  $\text{rank}(\rho) = r$ . Let  $P_{\pm} \in \mathbb{H}_d$  denote the orthogonal projectors onto the positive and negative parts of  $X$ . Then

$$\|X\|_1 = \langle P_+, X \rangle - \langle P_-, X \rangle = 2\langle P_+, X \rangle,$$

where the last equation follows from the fact that  $X$  is traceless. The matrix Hölder inequality asserts that  $\langle P_+, X \rangle \leq \|P_+\|_1 \|X\|$ . Therefore,

$$\|\rho - \sigma\|_1 = \|X\|_1 \leq 2\text{tr}(P_+) \|\rho - \sigma\| \leq 2r \|\rho - \sigma\|.$$

Indeed, the range of the orthogonal projector  $P_+$  has dimension at most  $r$ . The result follows once we establish that  $\|\rho - \sigma\| \leq 2\|\rho - M\|$ . This relation follows from the assumption that  $M$  has unit trace, but the proof is somewhat less transparent. ■

The following convergence bound is an immediate consequence of Proposition 3.8 and Lemma 3.9.

**Theorem 3.10 (Projected sample average estimator).** Suppose that we perform near-isotropic quantum measurements on identical copies of a quantum system that has the rank- $r$  density matrix  $\rho$ . Then the projected matrix sample average (3.7) obeys

$$\mathbb{P} \{ \|\hat{\rho}_n - \rho\|_1 \geq t \} \leq 2d \exp \left( \frac{-3nt^2}{256r^2d} \right)$$

The probability bound is valid for all  $t \geq 0$ .

In short, the matrix Bernstein inequality leads quickly to a strong error bound on the projected sample average estimator of a quantum state. The following observation is an immediate consequence of Theorem 3.10.

**Corollary 3.11** Fix a rank- $r$  density matrix  $\rho \in \mathcal{S}(\mathbb{H}_d)$ . Choose parameters  $\varepsilon, \delta \in (0, 1)$ . Then a total of

$$n \geq 86r^2 d \varepsilon^{-2} (\log(2d) + \log(1/\delta))$$

measurement repetitions (samples) are sufficient to guarantee that the projected sample average estimator obeys  $\|\hat{\rho}_n - \rho\|_1 \leq \varepsilon$  with probability at least  $1 - \delta$ .

We conclude that the projected sample average estimator almost saturates the fundamental lower bound (Theorem 3.7) on the sample complexity of any quantum state tomography procedure. Moreover, the performance is optimal up to a constant factor in the regime where the probability of success is at least  $1 - d^{-1}$ !



### 3.3.5 Generalization: Projected least squares

The matrix sample average estimator for near-isotropic quantum measurements is a special case of a general and practical procedure for quantum state tomography, called *projected least squares*. Here is a summary of this approach:

1. Fix a tomographically complete measurement  $\{\mathbf{H}_{\lambda_i} : 1 \leq i \leq m\}$ .
2. Estimate the probabilities  $p_i = \mathbb{P}\{\lambda_i | \boldsymbol{\rho}\}$  by *frequencies*. That is, prepare  $n$  identical realizations of the quantum system, measure them separately, and set

$$f_i^{(n)} = \frac{\# \text{ number of times outcome } \lambda_i \text{ was observed}}{\text{total number of measurements } n}.$$

3. Construct the least-squares estimator that results from replacing the true probabilities in Born's rule (3.1) by the frequency approximations:

$$\bar{\mathbf{R}}_n = \arg \min_{\mathbf{X} \in \mathbb{H}_d} \sum_{i=1}^m |f_i^{(n)} - \langle \mathbf{H}_{\lambda_i}, \mathbf{X} \rangle|^2.$$

4. Compute the Frobenius-norm projection of  $\bar{\mathbf{R}}_n$  onto the set  $\mathcal{S}(\mathbb{H}_d)$  of quantum states.

This procedure also results in a near-optimal quantum state estimator. As above, the analysis relies on the matrix Bernstein inequality. The main difference is that the solution to the linear inverse problem has a more complicated form when the measurement is not near-isotropic.





## 4. Graph Laplacians

"Pipes various," Wikimedia Commons

This lecture contains the fundamentals of spectral graph theory and harmonic analysis on graphs. The presentation is inspired by Dan Spielman's Fall 2018 course on spectral graph theory [[Spi](#)], Yuval Wigderson's notes on harmonic functions on graphs [[Wig](#)], and Rasmus Kyng's dissertation [[Kyn17](#)]. Any errors are my own.

A combinatorial graph encodes pairwise relationships among a family of objects. Graphs have intrinsic mathematical interest, as well as numerous computational applications. This lecture introduces the concept of a multigraph and the associated Laplacian matrix. The Laplacian encodes structural properties of the multigraph, and it can be understood with physical analogies to electrical networks.

Laplacian matrices play a role in learning methods based on harmonic analysis on manifolds. They also arise from the discretization of elliptic PDEs. The ultimate goal of this course is to present an efficient algorithm for solving a linear system in a graph Laplacian matrix, which can be used for both of the applications mentioned in this paragraph.

### 4.1 Multigraph basics

We will be working with (undirected) multigraphs, which are a lot like graphs, except that there may be many edges connecting a pair of vertices. This level of generality is important for us, so we must suffer the extra complexity.

#### 4.1.1 Undirected multigraphs

Let  $V$  be a set of  $n$  points, called *vertices*. The letters  $u$  and  $v$  will denote vertices. We may as well assume that  $V = \{1, \dots, n\}$ , which allows us place the vertices in order.

A *multiedge* is an unordered pair  $e = \{u, v\}$  of two distinct vertices  $u, v \in V$ . A multiedge represents an undirected link between the two vertices, and we forbid loops that connect a vertex to itself. It is convenient to abbreviate  $e = uv = vu$  for any multiedge connecting  $u$  and  $v$ . The notations  $u \in e$  and  $e \ni u$  both mean that the multiedge  $e$  contains the vertex  $u$ . We also say that  $e$  is *incident* on  $u$ .

We assign each multiedge a unique label so we can tell it apart from other multiedges between the same two vertices. At the risk of some confusion, we completely suppress this label from the notation.

An (undirected) *multigraph*  $G$  consists of a ground set  $V$  of vertices, along with a family  $E$  of multiedges. The letter  $m = |E|$  will refer to the total number of multiedges.

Somewhat abusively, we may write either  $e \in E$  or  $e \in G$  to indicate that the multigraph contains the multiedge  $e$ . (There are further notational abuses to come!)

We also equip with the multigraph  $G$  with a nonnegative weight function  $w : E \rightarrow \mathbb{R}_{++}$  that assigns a strictly positive value to each multiedge. Note that each multiedge joining a single pair of vertices can have a distinct weight.

We will always be working with the same ground set  $V$  of vertices, but there will be many multigraphs floating around. Therefore, it is often useful to qualify our notation by specifying a multigraph. For example, we may write  $w_G(e)$  or  $w(e, G)$  to refer to the weight of a multiedge in the multigraph  $G$ .

### 4.1.2 Connected components

A vertex  $u$  is a *neighbor* of a vertex  $v$  if the multigraph contains at least one multiedge  $e = uv$  linking the vertices  $u$  and  $v$ . We write  $u \sim v$  or  $v \sim u$  to indicate that  $u$  and  $v$  are neighbors.

We can iterate the neighbor relation to obtain *multi-hop neighborhoods* of a vertex. For a vertex  $u \in V$ , iteratively define

$$N^0(u) = \{u\} \quad \text{and} \quad N^k(u) = \{v' \in V : v' \sim v \text{ and } v \in N^{k-1}(u)\} \quad \text{for } k \in \mathbb{N}.$$

The set  $N^k(u)$  contains the vertices that are reachable from  $u$  by traversing exactly  $k$  multiedges. It is common to abbreviate  $N(u) = N^1(u)$ .

The *connected component*  $N^\infty(u)$  of a vertex  $u$  is the set of all vertices that are reachable from  $u$  via the neighbor relation:

$$N^\infty(u) = \bigcup_{k=0}^{\infty} N^k(u).$$

Every multigraph can be partitioned into a disjoint family of connected components. The relation  $N^\infty(u) = V$  means that every vertex in the multigraph is reachable from  $u$ . In the latter case, every vertex is reachable from every other vertex, and we say that the multigraph is *connected*.

From now on, we will assume that the multigraph  $G$  is connected.

### 4.1.3 Multidegree and total weight

The *degree*,  $\deg(u)$ , of a vertex  $u$  in the multigraph  $G$  is the total number of multiedges incident on  $u$ . That is,

$$\deg(u) = \deg(u, G) = |\{e \in G : e \ni u\}|.$$

Note that the multidegree of  $u$  need not coincide with the number of vertices that neighbor  $u$ .

The *total weight*  $w(u)$  of a vertex  $u$  in the multigraph  $G$  is the the sum of the weights of the multiedges that are incident on  $u$ . That is,

$$w(u) = w_G(u) = \sum_{e \in G, e \ni u} w_G(e).$$

Take care that the weight function has a different definition when applied to vertices and edges.

### 4.1.4 Interpretation: Plumbing

We can interpret a multigraph  $G = (V, E, w)$  as a plumbing network that connects the fixtures listed in  $V$  with the pipes listed in  $E$ . There may be many pipes connecting the same two fixtures. The weight  $w(e)$  associated with a pipe  $e$  increases with the “size” of the pipe.

For later reference, recall that the rate of flow along a pipe is proportional to the “size” of the pipe times the difference in pressure at the two endpoints. (The size of a circular pipe is the fourth power of the radius divided by the length.) This is called the Hagen–Poiseuille law.

### 4.1.5 Interpretation: Resistor networks

We can interpret a multigraph  $G = (V, E, w)$  as a wiring diagram that connects the terminals  $V$  with the wires  $E$ . There may be many wires connecting the same two terminals in parallel. The weight  $w(e)$  of a wire  $e$  is proportional to the electrical conductance of the wire. The weight  $w(e)$  is inversely proportional to the electrical resistance.

For later reference, we recall Ohm’s law:  $V = IR$ . In words, the difference ( $V$ ) in voltage at two terminals is proportional to the electrical current ( $I$ ) flowing between the terminals times the electrical resistance ( $R$ ) of the wire.

### 4.1.6 Example: A random walk

There is a natural construction of a random walk on  $G$ . Let  $u_0 \in V$  be the initial vertex. At each time  $k \in \mathbb{N} \cup \{0\}$ , we are at vertex  $u_k$ , and we draw the next vertex  $u_{k+1}$  in the walk at random according to the probability distribution

$$\mathbb{P}\{u_{k+1} = v \mid u_k = u\} = \frac{1}{w(u)} \sum_{e=uv \in G} w(e) \quad \text{for each } v \in N(u).$$

Each multiedge of the form  $e = uv$  appears once in the sum! In other words, we randomly choose one of the multiedges incident on  $u$  with probability in proportion to

its weight, and we traverse this edge to arrive at a new vertex  $v$ . This process repeats indefinitely.

The transition matrix  $\mathbf{Q}$  of the random walk is called the *random walk normalized Laplacian*, and it is obtained by diagonal reweighting of the ordinary Laplacian:

$$\mathbf{Q} = \text{diag}(w(u) : u \in \mathbf{V})^{-1} \mathbf{L}.$$

One can understand many features of the random walk by studying the eigenvalues and eigenvectors of the random walk normalized Laplacian. But this is a subject for another day.

## 4.2 Laplacian basics

Every multigraph is associated with a psd matrix, called the *Laplacian*. The properties of this matrix, as a linear operator, are intertwined with the structure of the multigraph.

### 4.2.1 The Laplacian of a multigraph

Let  $e = uv$  be a multiedge connecting distinct vertices  $u, v \in \mathbf{V}$ . The *elementary Laplacian* induced by the multiedge  $e$  is the matrix

$$\Delta_e = \Delta_{uv} = (\delta_u - \delta_v)(\delta_u - \delta_v)^* \in \mathbb{H}_{\mathbf{V}}.$$

Recall that  $\delta_u$  denotes the standard basis vector at vertex  $u$ . Observe that the elementary Laplacian is a psd matrix. In addition, the null space of the elementary Laplacian contains the constant vector  $\mathbf{1} \in \mathbb{R}^{\mathbf{V}}$ .

**Definition 4.1 (Graph Laplacian).** The *Laplacian* of the multigraph  $\mathbf{G}$  is the matrix

$$\mathbf{L} = \mathbf{L}_{\mathbf{G}} = \sum_{e \in \mathbf{G}} w(e) \Delta_e \in \mathbb{H}_{\mathbf{V}}. \quad (4.1)$$

The Laplacian  $\mathbf{L}$  is a psd matrix because it is a nonnegative sum of psd matrices. For distinct vertices  $u, v \in \mathbf{V}$ , the  $uv$  off-diagonal entry of the Laplacian records (the negative of) the total weight of all the multiedges connecting  $u$  and  $v$ :

$$(\mathbf{L})_{uv} = - \sum_{e=uv} w(e).$$

Meanwhile, the diagonal of the Laplacian records the total weight of each vertex of the graph:

$$w(u) = (\mathbf{L})_{uu} = \sum_{e \ni u} w(e) \quad \text{for each } u \in \mathbf{V}.$$

The diagonal and off-diagonal entries are related as

$$w(u) = (\mathbf{L})_{uu} = - \sum_{v \neq u} (\mathbf{L})_{uv}.$$

The last display is another statement of the fact that  $\mathbf{L}\mathbf{1} = \mathbf{0}$ .

**Exercise 4.1** Assume that  $\mathbf{G}$  is a connected multigraph. Prove that  $\text{null}(\mathbf{L}_{\mathbf{G}}) = \text{lin}\{\mathbf{1}\}$ .

**Exercise 4.2** Consider a symmetric matrix  $\mathbf{M} \in \mathbb{H}_{\mathbf{V}}$  for which

1.  $\mathbf{M}$  has nonnegative diagonal entries;
2.  $\mathbf{M}$  has nonpositive off-diagonal entries;
3.  $\mathbf{M}\mathbf{1} = \mathbf{0}$ .

Show that  $\mathbf{M}$  is the Laplacian of some (multi)graph. In particular, the class of Laplacian matrices forms a convex cone. (That is, the class is closed under addition and nonnegative scaling.)

#### 4.2.2 Correspondence between multigraphs and Laplacians

Each multigraph determines a unique Laplacian matrix, but the converse is not true. For the purposes of our presentation, we will elide this point by treating the Laplacian of the multigraph as a sum over multiedges. Moreover, we usually regard the multigraph and the Laplacian as interchangeable.

Let us take a minute to justify this decision more rigorously. We will construct a pair of matrices that are closely related to the Laplacian and that completely determine the multigraph. This approach is also useful for implementing algorithms.

To that end, let us enumerate the multiedges in the multigraph:  $e_1, e_2, \dots, e_m \in \mathbf{E}$ . The ordering is arbitrary, but fixed. The signed vertex–multiedge adjacency matrix  $\mathbf{A} \in \mathbb{R}^{V \times E}$  encodes the connectivity of the graph. The  $j$ th multiedge  $e_j = u_j v_j$  determines the  $j$ th column of the matrix:

$$\mathbf{a}_{:j} = \delta_{u_j} - \delta_{v_j} \quad \text{where} \quad u_j < v_j \quad \text{for each } e_j \in \mathbf{E}.$$

The ordering is chosen for concreteness, but it is unimportant. Second, introduce a nonnegative diagonal matrix  $\mathbf{W} \in \mathbb{H}_E$  that encodes the weights in the obvious way:

$$w_{jj} = w(e_j) \quad \text{for each } e_j \in \mathbf{E}.$$

Together,  $\mathbf{A}$  and  $\mathbf{W}$  contain all of the data about the graph.

These two matrices provide another construction of the Laplacian of the multigraph:

$$\mathbf{L} = \mathbf{A}\mathbf{W}\mathbf{A}^*.$$

This gives another precise sense to our identification of the Laplacian with a sum of multiedges.

#### 4.2.3 Projectors and pseudoinverses

We will be keenly interested in solving linear systems involving the Laplacian matrix  $\mathbf{L}$  of a multigraph  $\mathbf{G}$ . This requires some care because the Laplacian is singular.

**Definition 4.2 (Range projector and pseudoinverse).** Let  $\mathbf{G}$  be a connected multigraph with Laplacian  $\mathbf{L}$ . The orthogonal projector  $\mathbf{P} \in \mathbb{H}_V$  onto the range of  $\mathbf{L}$  is the matrix

$$\mathbf{P} = \mathbf{I} - |\mathbf{V}|^{-1} \mathbf{1}\mathbf{1}^*.$$

The *pseudoinverse*  $\mathbf{L}^\dagger \in \mathbb{H}_V$  is the unique psd matrix that satisfies

$$\mathbf{L}\mathbf{L}^\dagger = \mathbf{P} \quad \text{and} \quad \text{range}(\mathbf{L}^\dagger) = \text{range}(\mathbf{L}).$$



The next two results are easy consequences of the definitions.

**Exercise 4.3 (Laplacian pseudoinverse).** For a connected multigraph  $G$ , the Laplacian matrix  $L$  and the range projector  $P$  enjoy the following relationships:

1.  $LP = PL$ .
2.  $L^\dagger L = P$ .
3.  $L^\dagger LL^\dagger = L^\dagger$ .
4.  $LL^\dagger L = L$ .

**Exercise 4.4 (Laplacian linear systems).** Let  $G$  be a connected multigraph with Laplacian matrix  $L$ . Suppose that  $f \in \mathbb{R}^V$  satisfies the orthogonality relation  $\mathbf{1}^* f = 0$ . Then

$$Lx = f \quad \text{and} \quad \mathbf{1}^* x = 0 \quad \text{if and only if} \quad x = L^\dagger f.$$

#### 4.2.4 The Dirichlet form

The Laplacian induces a quadratic form, called the *Dirichlet form*:

$$\|x\|_L^2 = x^* L x = \sum_{e=uv} w(e) (x_u - x_v)^2 \quad \text{for } x \in \mathbb{R}^V.$$

Note that each multiedge of the form  $e = uv$  appears once in the sum! The associated pseudonorm is called the *Dirichlet energy*:

$$\|x\|_L = (x^* L x)^{1/2} \quad \text{for } x \in \mathbb{R}^V.$$

The Dirichlet norm of a vector  $x \in \mathbb{R}^V$  reflects its smoothness with respect to the graph structure.

The Dirichlet energy has various physical interpretations that are useful for constructing graph embeddings. The Dirichlet energy also provides a natural way to quantify the error in solving a linear system in the Laplacian matrix.

#### 4.2.5 Example: Laplacians and cuts

Here is a simple connection between the Dirichlet form and the combinatorial properties of a graph. A *cut* in a multigraph is a subset  $U$  of the vertices. The *weight* of a cut is the total weight of the multiedges that cross the cut:

$$\text{weight}(U) = \sum_{e=uv; u \in U; v \notin U} w(e).$$

Note that each multiedge  $e = uv$  in the multigraph appears at most once in the sum, with the orientation  $u \in U$  and  $v \notin U$ . The Laplacian allows us to express the weight of a cut. Evaluate the Dirichlet form at the indicator vector of the cut to obtain the weight of the cut:

$$\text{weight}(U) = \|x_U\|_L = x_U^* L x_U \quad \text{where} \quad x_U(u) = \begin{cases} 1, & u \in U; \\ 0, & u \notin U. \end{cases}$$

This formula allows us to use algebra to study combinatorial problems.

### 4.3 Harmonic analysis on multigraphs

We are now prepared to introduce the basic theory of harmonic functions on graphs.

#### 4.3.1 Harmonic functions

Harmonic functions arise as the solutions to homogeneous linear equations involving the Laplacian matrix.

**Definition 4.3 (Harmonic function).** Let  $G$  be a multigraph, and let  $U \subseteq V$  be a subset of the vertices. A function  $\varphi : V \rightarrow \mathbb{R}$  is *harmonic* on  $U$  if

$$(L\varphi)(u) = 0 \quad \text{for each } u \in U.$$

In particular, we say that the function  $\varphi$  is harmonic at a vertex  $u$  if  $(L\varphi)(u) = 0$ .

Our first result provides more intuition: A function  $\varphi$  is harmonic at a vertex  $u$  when the value  $\varphi(u)$  is the weighted average of the values  $\varphi(v)$  at its neighbors  $v \in N(u)$ .

**Proposition 4.4 (Averaging property).** The function  $\varphi : V \rightarrow \mathbb{R}$  is harmonic at a vertex  $u \in V$  if and only if

$$\varphi(u) = \frac{1}{w(u)} \sum_{e=uv} w(e) \varphi(v). \quad (4.2)$$

Each distinct multiedge of the form  $e = uv$  appears once in the sum!

*Proof.* This statement follows immediately from the definition (4.1) of the Laplacian and the definition (4.2) of harmonicity. ■

#### 4.3.2 Example: Hitting probabilities

Let  $B \subseteq V$  be a distinguished set of vertices. For a starting point  $u \in V$  and a vertex  $b \in B$ , the *hitting probability*  $h_b(u)$  is the probability that a random walk with initial vertex  $u_0 = u$  arrives at  $b$  before it arrives at any other vertex of  $B$ . Note that

$$h_b(b) = 1 \quad \text{and} \quad h_b(a) = 0 \quad \text{for each } a \in B \setminus \{b\}. \quad (4.3)$$

For each remaining vertex  $u \notin B$ , the hitting probability satisfies a simple recursion:

$$\begin{aligned} h_b(u) &= \sum_{v \in N(u)} \mathbb{P}\{u_1 = v \mid u_0 = u\} \cdot h_b(v) \\ &= \frac{1}{w(u)} \sum_{e=uv} w(e) h_b(v). \end{aligned} \quad (4.4)$$

Proposition 4.4 now implies that the hitting probability  $h_b$  is harmonic on  $V \setminus B$ .

#### 4.3.3 The maximum principle

The averaging property in Proposition 4.4 has a very significant consequence.

**Theorem 4.5 (Maximum principle).** Let  $G$  be a connected multigraph. If  $\varphi : V \rightarrow \mathbb{R}$  is harmonic on  $V$ , then  $\varphi$  is a constant function.

*Proof.* Suppose that  $\varphi$  is not constant. Introduce the set  $M$  of vertices where  $\varphi$  achieves its maximum value:

$$M = \arg \max \{ \varphi(u) : u \in V \}.$$

Since  $\varphi$  is not constant,  $M$  is a proper subset of  $V$ . Moreover, since  $G$  is connected, we can extract adjacent vertices  $u \sim u'$  where  $u \in M$  and  $u' \notin M$ . We calculate that

$$\sum_{e=uv} w(e) \varphi(v) < \left[ \sum_{e=uv} w(e) \right] \cdot \max_{v \in N(u)} \varphi(v) = w(u) \cdot \varphi(u).$$

Indeed, there is a multiedge  $uu'$  that participates in the sum, and  $\varphi(u') < \max \{ \varphi(v) : v \in N(u) \} = \varphi(u)$ . Equivalently,

$$\varphi(u) > \frac{1}{w(u)} \sum_{e=uv} w(e) \varphi(v).$$

Therefore,  $\varphi$  is not harmonic at  $u$ . We reject the hypothesis that  $\varphi$  is constant. ■

#### 4.3.4 Poles

Let us explain why Theorem 4.5 is called a maximum principle.

**Definition 4.6 (Pole).** Let  $\varphi : V \rightarrow \mathbb{R}$  be a function. A vertex  $v \in V$  is called a *pole* of the function if  $\varphi$  is *not* harmonic at  $v$ .

**Corollary 4.7 (Existence of poles).** Let  $G$  be a connected multigraph. If  $\varphi : V \rightarrow \mathbb{R}$  is a nonconstant function, then  $\varphi$  attains its maximum and minimum value at poles. In particular,  $\varphi$  has at least two poles.

*Proof.* In the proof of Theorem 4.5, we defined the set  $M$  of vertices where a function  $\varphi$  achieves its maximum value. We proved that  $M$  contains a vertex  $u$  where  $\varphi$  is not harmonic. Therefore, the function  $\varphi$  has a pole, and the maximum occurs there.

We can apply the same argument to the negation  $-\varphi$  to identify a pole  $u'$  where  $\varphi$  achieves its minimum.

Since  $\varphi$  is not constant, the maximum and minimum are not achieved at the same location. Thus  $u \neq u'$ . We conclude that  $\varphi$  has at least two poles. ■

#### 4.3.5 Harmonic extensions

Next, let us consider what happens if we require a harmonic function to meet some boundary conditions.

**Definition 4.8 (Harmonic extension).** Let  $\varphi_0 : B \rightarrow \mathbb{R}$  be a function on a nonempty set  $B \subseteq V$  of vertices. A *harmonic extension* of  $\varphi_0$  is a function  $\varphi : V \rightarrow \mathbb{R}$  that solves the linear system

$$\begin{cases} (L\varphi)(u) = 0, & u \in V \setminus B; \\ \varphi(u) = \varphi_0(u), & u \in B. \end{cases}$$

We can construct a unique harmonic extension under minimal hypotheses.

**Theorem 4.9 (Harmonic extensions).** Let  $G$  be a connected multigraph. Distinguish a nonempty set  $B \subseteq V$  of vertices. For any boundary data  $\varphi_0 : B \rightarrow \mathbb{R}$ , there is a unique harmonic extension of  $\varphi_0$  to the set  $V \setminus B$  of remaining vertices.

*Proof. Uniqueness:* Let  $\varphi_1$  and  $\varphi_2$  be two harmonic extensions of  $\varphi_0$ . Consider their difference  $\psi = \varphi_1 - \varphi_2$ . By linearity,  $\psi$  is harmonic on  $V \setminus B$ :

$$(L\psi)(u) = (L\varphi_1)(u) - (L\varphi_2)(u) = 0 \quad \text{for } u \in V \setminus B.$$

Corollary 4.7 implies that  $\psi$  achieves its maximum and minimum on  $B$ . But  $\psi$  has zero boundary data:

$$\psi(v) = \varphi_1(v) - \varphi_2(v) = 0 \quad \text{for each } v \in B.$$

Therefore,  $\psi$  is identically equal to zero.

**Existence:** Using the hitting probabilities (Section 4.3.2), we define the real-valued function

$$\varphi(u) = \sum_{b \in B} \varphi_0(b) h_b(u) \quad \text{for each } u \in V.$$

By the property (4.3) of the hitting probability  $h_b$ , the function  $\varphi$  agrees with  $\varphi_0$  on  $B$ . Meanwhile, for  $u \in V \setminus B$ , the recursion (4.4) gives

$$\begin{aligned} \varphi(u) &= \sum_{b \in B} \varphi_0(b) \frac{1}{w(u)} \sum_{e=uv} w(e) h_b(v) \\ &= \frac{1}{w(u)} \sum_{e=uv} w(e) \sum_{b \in B} \varphi_0(b) h_b(v) \\ &= \frac{1}{w(u)} \sum_{e=uv} w(e) \varphi(v). \end{aligned}$$

Therefore,  $\varphi$  is a harmonic extension of  $\varphi_0$  from  $B$  to  $V$ . ■

#### 4.3.6 Interpretation: Plumbing

Let  $p \in \mathbb{R}^V$  denote the pressure at each fixture in a network of pipes. Suppose that the network contains an inlet  $u_{\text{in}}$  where the pressure  $p_0(u_{\text{in}}) > 0$ , usually called a *source*. Suppose that the network also contains an outlet  $u_{\text{out}}$  where the pressure  $p_0(u_{\text{out}}) < 0$ , usually called a *sink*. The other fixtures are called *internal nodes*.

The theory of hydrodynamics states that the total (signed) flow  $f(u)$  at an internal node  $u$  equals zero because any water that enters must also leave. Each pipe incident on  $u$  contributes to the flow in or out of the fixture  $u$ . The rate of flow along a pipe  $e = uv$  is (proportional to) the size  $w(e)$  of the pipe times the difference  $p(u) - p(v)$  in pressure at the endpoints. Altogether,

$$0 = f(u) = \sum_{e=uv} w(e) (p(u) - p(v)) \quad \text{for each internal } u \in V.$$

We can rewrite this equation as

$$p(u) = \frac{1}{w(u)} \sum_{e=uv} w(e) p(v) \quad \text{for each internal } u \in V.$$

In other words, the pressure  $\mathbf{p}$  is harmonic at each internal node.

These equations can be combined:

$$\begin{cases} (\mathbf{L}\mathbf{p})(u) = 0, & u \text{ is internal;} \\ p(u) = p_0(u), & u \in \{u_{\text{in}}, u_{\text{out}}\}. \end{cases}$$

In summary, the pressure  $\mathbf{p} \in \mathbb{R}^V$  is the harmonic extension of the pressure at the source and sink. If there are many sources and sinks, a similar formula is valid.

#### 4.3.7 Interpretation: Resistor networks

Let  $\varphi \in \mathbb{R}^V$  denote the voltage at each node in an electrical network. Suppose that the network contains a source  $u_{\text{in}}$  where the voltage  $\varphi_0(u_{\text{in}}) > 0$ ; for example, a battery. Suppose that the network also contains a sink  $u_{\text{out}}$  where the voltage  $\varphi_0(u_{\text{out}}) < 0$ ; for example, the ground. The other fixtures are called *internal nodes*.

The theory of resistor networks states that the total current  $f(u)$  flowing through an internal node  $u$  equals zero because there is no input or output. Each wire incident on  $u$  contributes to the current flowing in or out of the node  $u$ . The amount of current flowing along a wire  $e = uv$  is proportional to the conductance  $w(e)$  and the difference  $\varphi(u) - \varphi(v)$  in voltage at the endpoints (i.e., the difference in electrical potential). Altogether,

$$0 = f(u) = \sum_{e=uv} w(e) (\varphi(u) - \varphi(v)) \quad \text{for each internal } u \in V.$$

We can rewrite this equation as

$$\varphi(u) = \frac{1}{w(u)} \sum_{e=uv} w(e) \varphi(v) \quad \text{for each internal } u \in V.$$

In other words, the voltage  $\varphi$  is harmonic at each internal node.

These equations can be combined:

$$\begin{cases} (\mathbf{L}\varphi)(u) = 0, & u \text{ is internal;} \\ \varphi(u) = \varphi_0(u), & u \in \{u_{\text{in}}, u_{\text{out}}\}. \end{cases}$$

In summary, the voltage  $\varphi \in \mathbb{R}^V$  is the harmonic extension of the voltages at the source and sink. If there are many sources and sinks, a similar formula remains valid.

Conversely, we can consider a vector  $\mathbf{f} \in \mathbb{R}^V$  of external currents. The value  $f(u)$  is the amount of current entering (or leaving) the network at vertex  $u$ . The network cannot hold current, so we must assume that  $\mathbf{1}^* \mathbf{f} = \mathbf{0}$ . That is, any current that enters must also leave. Then the induced voltages  $\varphi \in \mathbb{R}^V$  at each node satisfy

$$\varphi = \mathbf{L}^\dagger \mathbf{f}.$$

One can easily verify that  $\varphi$  is harmonic, except at the nodes  $u$  where  $f(u) \neq 0$ . In addition, the total induced voltage  $\mathbf{1}^* \varphi = 0$ , which reflects the fact that only voltage differences between terminals play a role in determining the flow.



## 5. Effective Resistance

“Old radio resistors,” Wikipedia

This lecture is based on Dan Spielman’s Fall 2018 course on spectral graph theory [Spi].

The parallel between harmonic analysis on graphs and electrical networks suggests further analogies. In this lecture, we explore several important concepts that arise from this perspective. We first discuss the notion of the effective resistance between two vertices in a graph. Then we introduce the leverage of an edge, which is a reflection of its importance in determining the graph structure. Using these concepts, we demonstrate that every graph Laplacian can be approximated strongly by the Laplacian of a sparse graph. We realize this approximation by nonuniform randomized sampling.

### 5.1 Resistance distance

We have introduced the machinery of harmonic functions so that we can understand properties of the pseudoinverse of a Laplacian. Of course, the pseudoinverse plays a role in the solution of linear systems. But it also has interesting physical interpretations related to the properties of the electrical network determined by the graph.

#### 5.1.1 Effective resistance

We begin with an important definition.

**Definition 5.1 (Effective resistance).** Let  $G$  be a connected multigraph on a vertex set  $V$  and with Laplacian matrix  $L$ . For vertices  $u, v \in V$ , not necessarily distinct, the *effective resistance*  $\varrho(u, v)$  between the vertices  $u$  and  $v$  is the nonnegative number

$$\varrho(u, v) = (\delta_u - \delta_v)^* L^\dagger (\delta_u - \delta_v).$$



As usual,  $\delta_u$  is a standard basis vector, and  $^\dagger$  denotes the pseudoinverse.

To understand why this quantity is called the effective resistance, note that

$$\boldsymbol{\varphi} = \mathbf{L}^\dagger(\delta_u - \delta_v) \in \mathbb{R}^V$$

is the vector of induced voltages if we inject one unit of current at vertex  $u$  and extract one unit of current at vertex  $v$ . Then

$$\varrho(u, v) = (\delta_u - \delta_v)^* \boldsymbol{\varphi} = \varphi(u) - \varphi(v).$$

In other words,  $\varrho(u, v)$  is the voltage difference between the vertices  $u$  and  $v$ , per unit of current. In other words, we can interpret it as the resistivity of the entire network against passing one unit of current from  $u$  to  $v$ .

In the hydraulic analogy, we can think about injecting a unit-rate flow at the inlet  $u$  and extracting it at the outlet  $v$ . The whole plumbing network behaves like a pipe that shunts the fluid between these two fixtures. The number  $\varrho(u, v)$  reflects the effective “size” of this compound pipe.

### 5.1.2 Effective resistance is a metric

A wonderful fact is that the effective resistance induces a metric on the vertex set of a multigraph. This result is an easy consequence of the maximum principle for harmonic functions. It will play a central role in the algorithm for solving Laplacian systems.

**Theorem 5.2 (Effective resistance is a metric).** Let  $\mathbf{G}$  be a connected multigraph on the vertex set  $V$ . The effective resistance  $\varrho$  determines a metric on the vertices. More precisely, for all vertices  $t, u, v \in V$ , it holds that

1.  $\varrho(u, v) = 0$  if and only if  $u = v$ .
2.  $\varrho(u, v) = \varrho(v, u)$ .
3.  $\varrho(t, v) \leq \varrho(t, u) + \varrho(u, v)$ .

*Proof.* Let  $\mathbf{L}$  be the Laplacian of the multigraph  $\mathbf{G}$ . The first two properties are easy exercises. For the triangle inequality, we define the functions

$$\begin{aligned} \boldsymbol{\varphi}_{tu} &= \mathbf{L}^\dagger(\delta_t - \delta_u), \quad \text{harmonic on } V \setminus \{t, u\}; \\ \boldsymbol{\varphi}_{uv} &= \mathbf{L}^\dagger(\delta_u - \delta_v), \quad \text{harmonic on } V \setminus \{u, v\}; \\ \boldsymbol{\varphi}_{tv} &= \mathbf{L}^\dagger(\delta_t - \delta_v). \end{aligned}$$

By linearity, these functions are related as  $\boldsymbol{\varphi}_{tv} = \boldsymbol{\varphi}_{tu} + \boldsymbol{\varphi}_{uv}$ . Taking the inner product of this identity with  $\delta_t - \delta_v$  gives

$$\varrho(t, v) = (\delta_t - \delta_v)^* \boldsymbol{\varphi}_{tv} = (\delta_t - \delta_v)^* \boldsymbol{\varphi}_{tu} + (\delta_t - \delta_v)^* \boldsymbol{\varphi}_{uv}.$$

It remains to bound the right-hand side in terms of the effective resistances  $\varrho(t, u)$  and  $\varrho(u, v)$ . We can accomplish this via the maximum principle.

To that end, we note the relation

$$\varphi_{tu}(t) - \varphi_{tu}(u) = \varrho(t, u) \geq 0.$$



By the maximum principle (Corollary 4.7), the harmonic function  $\varphi_{tu}$  takes its maximum value at the pole  $t$  and its minimum at the pole  $u$ . Thus,

$$\begin{aligned} (\delta_t - \delta_v)^* \varphi_{tu} &= \varphi_{tu}(t) - \varphi_{tu}(v) \\ &\leq \varphi_{tu}(t) - \varphi_{tu}(u) = (\delta_t - \delta_u)^* \varphi_{tu} = \varrho(t, u). \end{aligned}$$

Similarly,

$$\begin{aligned} (\delta_t - \delta_v)^* \varphi_{uv} &= \varphi_{uv}(t) - \varphi_{uv}(v) \\ &\leq \varphi_{uv}(u) - \varphi_{uv}(v) = (\delta_u - \delta_v)^* \varphi_{uv} = \varrho(u, v). \end{aligned}$$

The result follows when we sequence the last three displays.  $\blacksquare$

### 5.1.3 An alternative representation

There is another way of writing the effective resistance that will be useful for us. Let us introduce another piece of notation.

**Definition 5.3 (Normalizing map).** Let  $\mathbf{G}$  be a connected multigraph with Laplacian matrix  $\mathbf{L}$ . Define the *normalizing map*

$$\Phi(\mathbf{M}) = \Phi_{\mathbf{G}}(\mathbf{M}) = (\mathbf{L}^\dagger)^{1/2} \mathbf{M} (\mathbf{L}^\dagger)^{1/2} \quad \text{for } \mathbf{M} \in \mathbb{H}_V.$$

The exponent  $^{1/2}$  extracts the unique psd square root. The normalizing map  $\Phi$  is associated with the Laplacian of a particular multigraph  $\mathbf{G}$ , which will remain fixed throughout our discussion.

Let us note some properties of this map. First,  $\Phi(\mathbf{L}) = \mathbf{P}$ , where  $\mathbf{P}$  is the orthogonal projector onto  $\text{range}(\mathbf{L})$ . The function  $\Phi$  is an example of a positive linear map. Among many other properties,

$$\mathbf{M} \succeq \mathbf{0} \quad \text{implies} \quad \Phi(\mathbf{M}) \succeq \mathbf{0}.$$

See the book [Bha07] for an introduction to the theory of positive linear maps.

The normalizing map gives us another mechanism for expressing the effective resistance between two vertices. Indeed, the effective resistance is the spectral norm of the normalized elementary Laplacian of the unit edge connecting the two vertices.

**Proposition 5.4 (Effective resistance).** Let  $\mathbf{G}$  be a connected multigraph on the vertex set  $V$  and with normalizing map  $\Phi$ . For vertices  $u, v \in V$ ,

$$\varrho(u, v) = \|\Phi(\Delta_{uv})\|.$$

As always,  $\|\cdot\|$  is the spectral norm.

*Proof.* Since the effective resistance is nonnegative,

$$\begin{aligned} \varrho(u, v) &= \|(\delta_u - \delta_v)^* (\mathbf{L}^\dagger)^{1/2} (\mathbf{L}^\dagger)^{1/2} (\delta_u - \delta_v)\| \\ &= \|(\mathbf{L}^\dagger)^{1/2} (\delta_u - \delta_v) (\delta_u - \delta_v)^* (\mathbf{L}^\dagger)^{1/2}\| = \|(\mathbf{L}^\dagger)^{1/2} \Delta_{uv} (\mathbf{L}^\dagger)^{1/2}\|. \end{aligned}$$

We make the transition to the second line using the relation  $\|\mathbf{M}\mathbf{M}^*\| = \|\mathbf{M}^*\mathbf{M}\|$ . Identify the normalizing map to complete the argument.  $\blacksquare$

**Exercise 5.1** Prove that  $\varrho(u, v) = \text{tr } \Phi(\Delta_{uv})$ .

### 5.1.4 Leverage of a multiedge

We are now prepared to introduce a notion of the importance of a multiedge to the graph structure.

**Definition 5.5 (Leverage).** Let  $G$  be a connected multigraph on a vertex set  $V$  and with normalizing map  $\Phi$ . For each multiedge  $e = uv$  with weight  $w(e)$ , the *leverage* of the multiedge  $e$  is the quantity

$$\ell(e) = w(e) \varrho(u, v) = w(e) \|\Phi(\Delta_e)\|.$$

As usual,  $\varrho$  is the effective resistance induced by the multigraph  $G$ .

**Proposition 5.6 (Leverage of a multiedge).** Let  $G$  be a connected multigraph. For each multiedge  $e$  in the multigraph, the leverage  $\ell(e) \leq 1$ .

*Proof.* Introduce the Laplacian  $L$  of the multigraph:

$$L = \sum_{e \in G} w(e) \Delta_e.$$

Apply the normalizing map to the last display:

$$P = \Phi(L) = \sum_{e \in G} w(e) \Phi(\Delta_e). \quad (5.1)$$

Since  $\Phi$  is a positive linear map,

$$P \succcurlyeq w(e) \Phi(\Delta_e) \quad \text{for each } e \in G.$$

Taking the spectral norm, for each multiedge  $e = uv \in G$ , we have

$$1 \geq w(e) \|\Phi(\Delta_e)\| = w(e) \varrho(u, v) = \ell(e).$$

The last identity follows from Proposition 5.4. ■

The basic idea is that the effective resistance  $\varrho(u, v)$  measures how much voltage we need to push a unit of current from the node  $u$  to the node  $v$ . Meanwhile, the weight  $w(e)$  is proportional to the conductance of a wire that connects  $u$  and  $v$ . As a consequence of Ohm's law, the leverage  $\ell(e)$  measures the fraction of current that travels along the wire  $e$  if we push one unit of current from  $u$  to  $v$ .

If there is only one way to get from  $u$  to  $v$ , all the current must pass along the wire  $e = uv$ , and the leverage equals one. (Think of the edge connecting the two ends of a barbell graph.) Conversely, if there are many ways to get from  $u$  to  $v$ , some of the current may follow other routes, and the leverage of  $e = uv$  can be small. (Think of the edges in a complete graph.)

Similarly, if two vertices are wired in parallel, each of the multiedges will carry an equal proportion of the current between the vertices. By increasing the number of multiedges, we thereby decrease the leverage of each one to an equal part of the total.

**Exercise 5.2 (Total leverage).** Assume that  $G$  is a connected multigraph on  $n$  vertices. Prove that the total of all leverage scores is  $n - 1$ . That is,

$$\sum_{e \in G} \ell(e) = n - 1.$$

**Hint:** Use the identity (5.1) and Exercise 5.1.

## 5.2 Approximating a Laplacian by sampling

As an application of these ideas, we will prove that every Laplacian is well-approximated by the Laplacian of a sparse graph. We construct the sparse graph by randomly sampling edges according to their leverage. The result and high-level approach are due to Spielman & Srivastava [SS11]. The main tool in our analysis is the matrix Bernstein inequality, Theorem 1.13.

### 5.2.1 Spectral approximation

Let  $L$  be the Laplacian of a connected graph  $G$  on a set  $V$  of  $n$  vertices, with normalizing map  $\Phi$ . Let  $S$  be the Laplacian of another graph on the same vertex set  $V$ . We are interested in a very strong notion of approximation between these two Laplacians.

**Definition 5.7 (Spectral approximation).** For  $\varepsilon \in (0, 1)$ , we say that  $S$  is a  $\varepsilon$ -spectral approximation of  $L$  when

$$(1 - \varepsilon) L \preceq S \preceq (1 + \varepsilon) L.$$

If  $S$  is a spectral approximation of  $L$ , then the two Laplacians represent graphs with similar properties. Among other things,

1. The effective resistance between a pair of vertices in  $S$  is comparable to the effective resistance between the same pair of vertices in  $L$ .
2. The value of every graph cut in  $S$  is comparable to the value of the same cut in  $L$ .
3. The solution to the linear system  $Sx = f$  is not very different from the solution of  $Lx = f$ .

The last fact will be very important when we talk about how to solve Laplacian linear systems efficiently.

It is convenient to convert the spectral approximation condition into another form that is more amenable to analysis.

**Proposition 5.8 (Spectral approximation).** Let  $L$  be the Laplacian of a connected graph, and let  $S$  be the Laplacian of another graph on the same vertex set. Suppose that

$$\|\Phi(S - L)\| \leq \varepsilon.$$

Then  $S$  is a spectral approximation of  $L$  with quality  $\varepsilon$ .

*Proof.* First, subtract  $L$  from both sides of the relation:

$$-\varepsilon L \preceq S - L \preceq +\varepsilon L.$$

Apply the normalizing map to this relation to obtain the equivalent condition

$$-\varepsilon P \preceq \Phi(S - L) \preceq +\varepsilon P.$$

Since the range of  $\Phi(S - L)$  is contained in the range of  $P$ , it is sufficient to prove that

$$\|\Phi(S - L)\| \leq \varepsilon \|P\| = \varepsilon.$$

The last step uses the fact that  $P$  is an orthogonal projector. ■

### 5.2.2 The sampling model

The representation of the graph Laplacian as a sum of weighted edges allows us to construct a matrix approximation by random sampling, similar to what we did in Lecture 2.

Recall that the Laplacian  $\mathbf{L}$  admits the representation

$$\mathbf{L} = \sum_{e \in \mathbf{G}} w(e) \Delta_e.$$

To construct a sparse Laplacian that approximates  $\mathbf{L}$ , we introduce a random elementary Laplacian:

$$\mathbf{R} = \frac{w(e)}{p_e} \Delta_e \quad \text{with probability } p_e > 0 \text{ for each } e \in \mathbf{G}.$$

It is immediate that  $\mathbb{E} \mathbf{R} = \mathbf{L}$ . For a parameter  $K \geq 1$ , we construct the Laplacian approximation by averaging  $K$  copies of  $\mathbf{R}$ :

$$\mathbf{S} = \frac{1}{K} \sum_{i=1}^K \mathbf{R}_i \quad \text{where } \mathbf{R}_i \sim \mathbf{R} \text{ iid.}$$

Then  $\mathbf{S}$  is the Laplacian of a weighted graph with at most  $K$  edges, and  $\mathbf{S}$  is an unbiased estimator of  $\mathbf{L}$ .

### 5.2.3 The sampling probabilities

Our goal is to control the size of  $\Phi(\mathbf{S} - \mathbf{L})$ , and we will exploit our analysis of matrix sampling estimators, Theorem 2.1. To activate this result, we need to make each summand uniformly bounded. To that end, calculate that

$$\|\Phi(\mathbf{R})\| = \frac{w(e)}{p_e} \|\Phi(\Delta_e)\| = \frac{\ell(e)}{p_e}.$$

Therefore, it is natural to select the sampling probabilities proportional to the leverage of the edges:  $p_e = c\ell(e)$ . The constant is selected to obtain a probability mass:

$$1 = \sum_e p_e = c \sum_e \ell(e) = (n-1)c.$$

We have used Exercise 5.2 here. On other words,  $c = 1/(n-1)$ .

### 5.2.4 The analysis

The analysis is easy now. The upper bound parameter in Theorem 2.1 satisfies

$$B = \sup \|\Phi(\mathbf{R})\| = \sup \frac{\ell(e)}{\ell(e)/(n-1)} = n-1.$$

The per-sample second moment satisfies

$$\begin{aligned} m_2(\Phi(\mathbf{R})) &= \|\mathbb{E} \Phi(\mathbf{R})^2\| \leq \|\mathbb{E} [\|\Phi(\mathbf{R})\| \cdot \Phi(\mathbf{R})]\| \\ &\leq (n-1) \|\Phi(\mathbb{E} \mathbf{R})\| = (n-1) \|\mathbf{P}\| = n-1. \end{aligned}$$

Theorem 2.1 immediately implies that

$$\begin{aligned} \mathbb{E} \|\Phi(\mathbf{S} - \mathbf{L})\| &\leq \sqrt{\frac{2m_2(\mathbf{R}) \log(2n)}{K}} + \frac{2B \log(2n)}{3K} \\ &\leq \sqrt{\frac{2(n-1) \log(2n)}{K}} + \frac{2(n-1) \log(2n)}{3K}. \end{aligned}$$

Set  $K = 4\varepsilon^{-2}(n-1) \log(2n)$  to arrive at the bound

$$\mathbb{E} \|\Phi(\mathbf{S} - \mathbf{L})\| < \varepsilon.$$

The final estimate assumes that  $\varepsilon \leq 1$ .

By the probabilistic method, every graph with sufficiently small leverage scores admits an  $\varepsilon$ -spectral approximation with at most  $4\varepsilon^{-2}n \log n$  edges. Note that this bound is *independent* of the number  $m$  of edges in the target graph  $\mathbf{G}$ !

Let us remark that the approach here tracks the presentation in Lecture 2. We recognize that the representation of the Laplacian as a sum of elementary Laplacians furnishes a breakdown of the matrix into simple components. Once we agree that our goal is to obtain an  $\varepsilon$ -spectral approximation, the matrix sampling result, Theorem 2.1, tells us exactly what properties the sampling probability ought to have. The leverage emerges as a formal consequence.

### 5.2.5 Computational aspects

This argument gives an algorithm for constructing a sparse graph that is a spectral approximation of an arbitrary graph. Unfortunately, to implement this method, we need to compute the leverages so that we can perform the sampling. The naïve approach to this problem involves  $\Theta(n^3)$  computation. The literature on theoretical algorithms contains techniques that can achieve this goal more efficiently, but these methods may not be entirely practical. See [Spi] for discussion and references.

### 5.2.6 Conclusion

To conclude, we have established the following theorem.

**Theorem 5.9 (Spielman & Srivastava, 2011).** Let  $\mathbf{G}$  be a connected graph on  $n$  vertices with Laplacian  $\mathbf{L}$ . Fix a parameter  $\varepsilon \in (0, 1)$ . Then there is a connected graph on the same vertex set, with at most  $4\varepsilon^{-2}n \log n$  edges, and whose Laplacian is an  $\varepsilon$ -spectral approximation to  $\mathbf{L}$ .

The analysis of this sampling estimator cannot be improved in general beyond the constants. Indeed, for a complete graph, the leverage of each edge is the same. The sampling technique chooses each edge with equal probability, so an individual sample is equally likely to be incident on each vertex. But the coupon collector problem tells us that we need  $\Theta(n \log n)$  samples to acquire edges incident on all  $n$  vertices. This outcome is prerequisite for  $\mathbf{S}$  to be the Laplacian of a connected graph, which is necessary if  $\mathbf{S}$  is to approximate  $\mathbf{L}$  spectrally.

Nevertheless, the sparsity bound in the theorem is not optimal for graph approximations. See [BSS14] for a sharp result based on a deterministic construction.



A. Cholesky

# Sur la résolution numérique des Systèmes d'équations linéaires

La solution des problèmes dépendant de données expérimentales, qui peuvent dans certains cas être soumises à des conditions, et auxquelles on applique la méthode des moindres carrés est toujours subordonnée au calcul numérique des racines d'un système d'équations linéaires. C'est le cas de la recherche des lois physiques ; c'est aussi le cas de la compensation de

## 6. Solving Laplacian Systems

This lecture is adapted from Rasmus Kyng's dissertation [Kyn17].

In this lecture, we will discuss computational methods for solving Laplacian linear systems. The classic direct method is based on computing the Cholesky decomposition of the Laplacian matrix. This decomposition takes a special form for the Laplacian, as compared with a general psd matrix. This special form will serve as the foundation for developing a very efficient algorithm for solving graph Laplacian systems, as we will see in Lecture 8.

### 6.1 Cholesky meets Laplace

This section gives an overview of a classic approach for solving a Laplacian system.

#### 6.1.1 Setup

Let  $G$  be a connected multigraph. The vertex set  $V = \{1, \dots, n\}$ . The multiedge set  $E$  comprises  $m$  edges. The weight function  $w_G : E \rightarrow \mathbb{R}_{++}$ . As always,  $L$  denotes the weighted Laplacian matrix of the graph. We will treat the multigraph and the Laplacian as interchangeable by presenting the Laplacian as a weighted sum of multiedges.

#### 6.1.2 Laplacian systems

Suppose that we are given a forcing vector  $f \in \mathbb{R}^V$  that is orthogonal to the constant vector  $\mathbf{1}^* f = 0$ . Our aim is to solve the linear system

$$Lx = f. \quad (6.1)$$

Write  $x_\star \in \mathbb{R}^V$  for the (unique) solution to this system with  $\mathbf{1}^* x_\star = 0$ .



### 6.1.3 Solution via Cholesky decomposition

A standard approach to solving a Laplacian linear system is to extract a Cholesky decomposition of the Laplacian:

$$\mathbf{L} = \mathbf{C}\mathbf{C}^* \quad \text{where } \mathbf{C} \text{ is lower-triangular.}$$

We will spend most of this lecture going over the details of how Cholesky decomposition works for Laplacian matrices. The cost of producing this decomposition is usually  $\Theta(n^3)$  arithmetic operations. Owing to fill-in, the lower-triangular factor  $\mathbf{C}$  often has fully  $\Theta(n^2)$  nonzero entries.

Once we have the triangular factorization, we can solve the linear system (6.1) with  $\Theta(n^2)$  arithmetic operations. Indeed,

$$\mathbf{x}_\star = \mathbf{L}^\dagger \mathbf{f} = (\mathbf{C}^*)^\dagger (\mathbf{C}^\dagger \mathbf{f}).$$

We can apply  $\mathbf{C}^\dagger$  using forward substitution in time  $\Theta(n^2)$ , and we can apply  $(\mathbf{C}^*)^\dagger$  using backward substitution in time  $\Theta(n^2)$ . This approach produces results that are accurate (almost) to machine precision [Higo2, Chaps. 8, 10].

## 6.2 Cholesky decomposition: Matrix view

Let us explain the process of computing the Cholesky decomposition of a psd matrix. We begin with a linear-algebraic treatment that is applicable to any matrix. In the next section, we specialize to the case of Laplacian matrices.

### 6.2.1 Setup

Let  $\mathbf{M} \in \mathbb{H}_n$  be a psd matrix. The Cholesky decomposition iteratively reduces the psd matrix to a product of lower-triangular factors.

### 6.2.2 First step of the Cholesky decomposition

Let us begin with a visual description of the first step in the Cholesky process. Writing out the first row and column explicitly, we can express the matrix  $\mathbf{M}$  as

$$\mathbf{M} = \begin{bmatrix} d & -\mathbf{a}^* \\ -\mathbf{a} & \mathbf{M}_2 \end{bmatrix}.$$

In this expression,  $d$  is a nonnegative number (because  $\mathbf{M}$  is psd),  $\mathbf{a} \in \mathbb{R}^{n-1}$  and  $\mathbf{M}_2 \in \mathbb{H}_{n-1}$ . Construct the rank-one psd matrix

$$d^\dagger \begin{bmatrix} d \\ -\mathbf{a} \end{bmatrix} \begin{bmatrix} d \\ -\mathbf{a} \end{bmatrix}^* = \begin{bmatrix} d & -\mathbf{a}^* \\ -\mathbf{a} & d^\dagger \mathbf{a} \mathbf{a}^* \end{bmatrix}.$$

(Since  $\mathbf{M}$  is psd, if the diagonal entry  $d = 0$ , then also  $\mathbf{a} = \mathbf{0}$ .) Therefore, we can eliminate the first row and column of  $\mathbf{M}$  by subtracting this rank-one matrix:

$$\mathbf{M}/1 = \mathbf{M} - d^\dagger \begin{bmatrix} d \\ -\mathbf{a} \end{bmatrix} \begin{bmatrix} d \\ -\mathbf{a} \end{bmatrix}^* = \begin{bmatrix} 0 & \mathbf{0}^* \\ 0 & \mathbf{M}_2 - d^\dagger \mathbf{a} \mathbf{a}^* \end{bmatrix}.$$

The notation  $\mathbf{M}/1$  refers to the Schur complement of  $\mathbf{M}$  with respect to the first coordinate subspace.

Observe that the reduced matrix  $\mathbf{M}/1$  remains psd. Indeed, for each vector  $\mathbf{x} \in \mathbb{R}^{n-1}$ , define  $\alpha = d^\dagger \mathbf{a}^* \mathbf{x}$ , and calculate that

$$0 \leq \begin{bmatrix} \alpha \\ \mathbf{x} \end{bmatrix}^* \mathbf{M} \begin{bmatrix} \alpha \\ \mathbf{x} \end{bmatrix} = d\alpha^2 - 2\alpha \mathbf{a}^* \mathbf{x} + \mathbf{x}^* \mathbf{M}_2 \mathbf{x} = \mathbf{x}^* (\mathbf{M}_2 - d^\dagger \mathbf{a} \mathbf{a}^*) \mathbf{x}.$$

Since  $\mathbf{x}$  is arbitrary, the Schur complement  $\mathbf{M}/1$  is psd.

We can apply this elimination procedure to each remaining coordinate in sequence, reducing the size of the nonzero block at each step.

### 6.2.3 Cholesky decomposition, without pivoting

Here is a more formal description of the Cholesky decomposition, where we eliminate the coordinates in lexicographic order.

To begin the process, set  $\mathbf{S}_0 = \mathbf{M}$ . At each step  $i = 1, 2, \dots, n$ , we eliminate the  $i$ th coordinate. We write  $u_i = i$  to emphasize the difference between the choice of coordinate ( $u_i$ ) and the step ( $i$ ) in the iteration procedure. Introduce the vector

$$\mathbf{c}_i = \frac{1}{\sqrt{(\mathbf{S}_{i-1})_{u_i u_i}}} \cdot \mathbf{S}_{i-1} \delta_{u_i}.$$

(If the diagonal entry  $(\mathbf{S}_{i-1})_{u_i u_i}$  happens to equal zero, we set  $\mathbf{c}_i = \mathbf{0}$ .) Zero out the row and column in  $\mathbf{S}_{i-1}$  indexed by the coordinate  $u_i$  by forming the Schur complement:

$$\mathbf{S}_i = \mathbf{S}_{i-1}/u_i = \mathbf{S}_{i-1} - \mathbf{c}_i \mathbf{c}_i^*.$$

We continue this process for  $n$  steps.

At each iteration  $i$ , the matrix  $\mathbf{S}_{i-1}$  is psd, and it has the block form

$$\mathbf{S}_{i-1} = \begin{bmatrix} \mathbf{0}_{i \times i} & \mathbf{0}_{i \times (n-i)} \\ \mathbf{0}_{(n-i) \times i} & \star \end{bmatrix}.$$

The symbol  $\star$  indicates an  $(n-i) \times (n-i)$  block of nonzero coordinates. In particular, after  $n$  steps of this procedure,  $\mathbf{S}_n = \mathbf{0}$ .

When the iteration is complete, we collect the vectors  $\mathbf{c}_i$  into a matrix:

$$\mathbf{C} = [\mathbf{c}_1 \quad \mathbf{c}_2 \quad \dots \quad \mathbf{c}_n] \in \mathbb{R}^{n \times n}.$$

Since  $\mathbf{S}_{i-1}$  is supported on the coordinates  $\{i, \dots, n\}$ , so is the vector  $\mathbf{c}_i$ . Therefore, the matrix  $\mathbf{C}$  is lower-triangular.

To understand the role of the matrix  $\mathbf{C}$ , observe that

$$\mathbf{M} = \mathbf{S}_0 - \mathbf{S}_n = \sum_{i=1}^n (\mathbf{S}_{i-1} - \mathbf{S}_i) = \sum_{i=1}^n \mathbf{c}_i \mathbf{c}_i^* = \mathbf{C} \mathbf{C}^*.$$

The factorization  $\mathbf{M} = \mathbf{C} \mathbf{C}^*$  is called a Cholesky decomposition of the input matrix  $\mathbf{M}$ . The Cholesky procedure progressively decomposes the input matrix:

$$\mathbf{M} = \mathbf{S}_i + \sum_{k=1}^i \mathbf{c}_k \mathbf{c}_k^* \quad \text{for each } i = 0, 1, 2, \dots, n.$$

This relation breaks down the process into the part of the matrix that remains to be factored ( $\mathbf{S}_i$ ) and the part of the factorization that is done (the sum of rank-one terms).

### 6.2.4 Cholesky decomposition, with pivoting

It is not necessary to eliminate the coordinates in the lexicographic order. At iteration  $i$ , suppose instead that we eliminate the coordinate  $u_i = \pi(i)$ , where  $\pi : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$  is a permutation (i.e., a bijection). The permutation  $\pi$  can be chosen in advance, or  $\pi(i)$  may be selected at iteration  $i$ . Otherwise, the decomposition algorithm is the same as before.

In this case, the computed matrix  $\mathbf{C}$  is morally lower-triangular. The permutation  $\pi$  gives the order of elimination for solving the system  $\mathbf{C}\mathbf{x} = \mathbf{f}$  by substitution. We omit further details.

### 6.2.5 Computational cost

At the  $i$ th step of the Cholesky decomposition, the cost of computing the Schur complement is  $\Theta((n - i)^2)$  arithmetic operations. Therefore, the total cost of  $n$  iterations is  $\Theta(n^3)$ .

When the input matrix is sparse, we may be able to economize during the early iterations by exploiting sparsity. Nevertheless, each time we form a Schur complement, the matrix often becomes denser, a process known as *fill-in*. It is hard to avoid fill-in, except in special cases. As a consequence, we generally need  $\Theta(n^3)$  operations to obtain the Cholesky decomposition. This is very expensive.

## 6.3 Cholesky decomposition: Graph view

When we apply the Cholesky algorithm to a Laplacian, we can interpret the basic step as a combinatorial operation on a multigraph.

### 6.3.1 Setup

Let  $\mathbf{G}$  be a connected multigraph with Laplacian matrix  $\mathbf{L}$ . We can interpret the Cholesky decomposition of the Laplacian matrix  $\mathbf{L}$  in graph-theoretic terms. In particular, we can express the matrices that arise during the process as Laplacians!

### 6.3.2 First step of the Cholesky decomposition

To illustrate the idea, suppose that we want to eliminate the first vertex. Isolating the role of the first vertex,

$$\mathbf{L} = \sum_{e \in \mathbf{G}} w(e) \Delta_e = \begin{bmatrix} d & -\mathbf{a}^* \\ -\mathbf{a} & \mathbf{L}_2 \end{bmatrix}.$$

Since  $\mathbf{L}$  is a Laplacian, we can say more about the terms that appear here:

$$d = w_{\mathbf{L}}(1) \geq 0 \quad \text{and} \quad \mathbf{a} \geq \mathbf{0} \quad \text{and} \quad \mathbf{a}^* \mathbf{1} = d.$$

Indeed, the first diagonal entry of the Laplacian is the total weight  $w_{\mathbf{L}}(1)$  of the first vertex. We regard  $\mathbf{a} \in \mathbb{R}^{\mathbf{V} \setminus \{1\}}$ . The off-diagonal entries in the Laplacian are nonpositive, and the number  $a_v$  is the total weight of all multiedges of the form  $e = 1v$  for each vertex  $v \neq 1$ . The identity  $\mathbf{a}^* \mathbf{1} = d$  reflects the fact that the diagonal entry  $w_{\mathbf{L}}(1)$  is

the sum of all the weights of multiedges incident on the first vertex. We write  $\mathbf{L}_2$  as a placeholder for the submatrix indexed by the vertices  $V \setminus \{1\}$ .

To compute the Schur complement  $\mathbf{L}/1$  with respect to the first vertex, we subtract a rank-one matrix from the Laplacian. Introduce the vector

$$\mathbf{c} = \frac{1}{\sqrt{d}} \mathbf{L} \delta_1 = \frac{1}{\sqrt{d}} \begin{bmatrix} d \\ -\mathbf{a} \end{bmatrix}, \quad \text{so that} \quad \mathbf{c} \mathbf{c}^* = \begin{bmatrix} d & -\mathbf{a}^* \\ -\mathbf{a} & d^\dagger \mathbf{a} \mathbf{a}^* \end{bmatrix}.$$

(If  $d = 0$ , we interpret the fraction as computing a pseudoinverse.) Then the Schur complement takes the form

$$\mathbf{L}/1 = \mathbf{L} - \mathbf{c} \mathbf{c}^* = \begin{bmatrix} 0 & \mathbf{0}^* \\ \mathbf{0} & \mathbf{L}_2 - d^\dagger \mathbf{a} \mathbf{a}^* \end{bmatrix}.$$

To understand what is really happening here, we will think about this operation as the composition of two steps.

Define the Laplacian of the set of multiedges incident on the first vertex:

$$\text{STAR}(1) = \sum_{e=1v} w(e) \Delta_e = \begin{bmatrix} d & -\mathbf{a}^* \\ -\mathbf{a} & \text{diag}(\mathbf{a}) \end{bmatrix}.$$

Here,  $\text{diag}(\mathbf{a})$  is the diagonal matrix determined by the vector  $\mathbf{a}$ . This Laplacian is called the *star* induced by the first vertex. Adding and subtracting  $\text{STAR}(1)$  from the Schur complement  $\mathbf{L}/1$ , we obtain

$$\mathbf{L}/1 = (\mathbf{L} - \text{STAR}(1)) + (\text{STAR}(1) - \mathbf{c} \mathbf{c}^*).$$

We will check that each of the large parentheses defines a Laplacian matrix. Since the class of Laplacians is closed under addition, the Schur complement  $\mathbf{L}/1$  is also a Laplacian matrix!

The first parenthesis is simply the Laplacian of the multigraph obtained by removing from  $G$  the multiedges incident on the first vertex:

$$\mathbf{L} - \text{STAR}(1) = \sum_{e \neq 1} w(e) \Delta_e.$$

This point follows immediately from the definitions. Observe that none of the remaining multiedges is incident on the first vertex. This is equivalent to the matrix  $\mathbf{L} - \text{STAR}(1)$  being supported on the coordinates  $v \neq 1$ .

Consider the second parenthesis:

$$\text{STAR}(1) - \mathbf{c} \mathbf{c}^* = \begin{bmatrix} 0 & \mathbf{0}^* \\ \mathbf{0} & \text{diag}(\mathbf{a}) - d^\dagger \mathbf{a} \mathbf{a}^* \end{bmatrix}.$$

This matrix is also a Laplacian! Indeed, by direct calculation, the diagonal entries are nonnegative, the off-diagonal entries are nonpositive, and each row sums to one. Alternatively, we can write

$$\begin{aligned} \begin{bmatrix} 0 & \mathbf{0}^* \\ \mathbf{0} & \text{diag}(\mathbf{a}) - d^\dagger \mathbf{a} \mathbf{a}^* \end{bmatrix} &= \frac{1}{2d} \sum_{v_1, v_2 \neq 1} a_{v_1} a_{v_2} (\delta_{v_1} - \delta_{v_2})(\delta_{v_1} - \delta_{v_2})^* \\ &= \frac{1}{2w_L(1)} \sum_{\substack{e_1=1v_1 \\ e_2=1v_2}} w(e_1) w(e_2) \Delta_{v_1 v_2}. \end{aligned}$$

(We interpret the fraction bar as computing a pseudoinverse.) This Laplacian is also called the *clique* induced by eliminating the vertex 1.

To wit, the process of computing the Schur complement of a Laplacian with respect to a vertex amounts to removing the star induced by the vertex and adding back the clique induced by eliminating the vertex.

### 6.3.3 Stars and cliques

Let us develop this construction in more generality. Let  $\mathbf{S}$  be the Laplacian of a multigraph on the vertex set  $V$ , expressed as a weighted sum of multiedges:

$$\mathbf{S} = \sum_{e \in \mathcal{S}} w_{\mathbf{S}}(e) \Delta_e \in \mathbb{H}_V.$$

It is rather irritating, but necessary, to keep track of which Laplacian we are operating on. The notation will reflect the choice.

Suppose that we wish to eliminate the vertex  $u$  from the Laplacian. Define the *star* induced by a vertex  $u$  is the Laplacian generated by the weighted edges in  $\mathbf{S}$  that are incident on  $u$ . That is,

$$\text{STAR}(u, \mathbf{S}) = \sum_{e=uv \in \mathcal{S}} w_{\mathbf{S}}(e) \Delta_e. \quad (6.2)$$

The sum takes place over all multiedges  $e$  in  $\mathbf{S}$  that are incident on the vertex  $u$ . The *clique* induced by eliminating the vertex  $u$  from  $\mathbf{S}$  is a weighted Laplacian

$$\text{CLIQUE}(u, \mathbf{S}) = \frac{1}{2w_{\mathbf{S}}(u)} \sum_{e_1=uv_1 \in \mathcal{S}} \sum_{e_2=uv_2 \in \mathcal{S}} w_{\mathbf{S}}(e_1) w_{\mathbf{S}}(e_2) \Delta_{v_1 v_2}. \quad (6.3)$$

Each sum takes place over all multiedges  $e$  in  $\mathbf{S}$  that are incident on the vertex  $u$ .

The star is the Laplacian of a multigraph; the clique is also the Laplacian of a multigraph. By a direct calculation, the star and the clique satisfy the identity

$$\text{CLIQUE}(u, \mathbf{S}) - \text{STAR}(u, \mathbf{S}) = -\frac{1}{w_{\mathbf{S}}(u)} (\mathbf{S}\delta_u)(\mathbf{S}\delta_u)^*. \quad (6.4)$$

Therefore, the Schur complement  $\mathbf{S}/u$  takes the form

$$\begin{aligned} \mathbf{S}/u &= \mathbf{S} - \frac{1}{w_{\mathbf{S}}(u)} (\mathbf{S}\delta_u)(\mathbf{S}\delta_u)^* \\ &= (\mathbf{S} - \text{STAR}(u, \mathbf{S})) + \text{CLIQUE}(u, \mathbf{S}). \end{aligned}$$

As before, the Schur complement  $\mathbf{S}/u$  is the Laplacian of a multigraph. This multigraph has no edges incident on the vertex  $u$ . Moreover, if a vertex does not participate in  $\mathbf{S}$ , it does not participate in  $\mathbf{S}/u$ .

To repeat, we compute the Schur complement of  $\mathbf{S}$  with respect to a vertex  $u$  by adding the clique induced by eliminating the vertex  $u$  from  $\mathbf{S}$  and then removing the star induced by the vertex  $u$  from  $\mathbf{S}$ .

### 6.3.4 Cholesky decomposition of a Laplacian

We are now prepared to summarize the process of computing the (pivoted) Cholesky decomposition of the Laplacian  $L$  of the multigraph  $G$ .

Set  $S_0 = L$ . For each iteration  $i = 1, 2, 3, \dots, n$ , select a vertex  $u_i$ . Extract the associated normalized column of the Laplacian:

$$\mathbf{c}_i = \frac{1}{\sqrt{(S_{i-1})_{u_i u_i}}} S_{i-1} \delta_{u_i}.$$

Compute the Schur complement with respect to the vertex  $u_i$ :

$$S_i = S_{i-1} / u_i = (S_{i-1} - \text{STAR}(u_i, S_{i-1})) + \text{CLIQUE}(u_i, S_{i-1}).$$

Let us emphasize that  $S_i$  remains a Laplacian matrix, but it has no multiedges incident on the vertices  $\pi(1), \dots, \pi(i)$ . We have reduced the size of the problem, and the process continues.

After  $n$  steps, construct the morally lower-triangular matrix

$$C = [\mathbf{c}_1 \quad \dots \quad \mathbf{c}_n] \in \mathbb{R}^{V \times V}.$$

Then the initial Laplacian admits the factorization

$$L = CC^*.$$

Last, we record the permutation  $\pi$  defined by  $\pi(i) = u_i$  for  $i = 1, \dots, n$ . This permutation reflects the order in which the vertices were eliminated, and it is also determines the substitution order for solving the linear system  $C\mathbf{x} = \mathbf{f}$ .

Let us remark that there is a standard approach to selecting vertices to eliminate from a graph Laplacian. At each step, we choose the remaining vertex that has the minimum degree.

### 6.3.5 An opportunity

Recall that the Cholesky decomposition is expensive because of the cost of computing the Schur complement. For general psd matrices, we compute the Schur complement by subtracting a rank-one matrix. It is not clear how to approximate this operation accurately.

For graph Laplacians, however, we expressed the Schur complement as the composition of two simple graph operations. Removing the star induced by a vertex is straightforward and inexpensive. The dominant cost arises from introducing the clique; this operation is quadratic in the number of edges incident on the vertex we eliminate.

Nevertheless, the clique is expressed as a weighted sum of many elementary Laplacians. As a consequence, we can try to approximate the clique by sampling. This is the core idea behind the `SparseCholesky` algorithm, which we detail in the Lecture 8.

The `SparseCholesky` is an iterative algorithm that constructs a sequence of random matrices. To analyze this kind of algorithm, we need more sophisticated matrix concentration tools. The next lecture turns to the subject of matrix martingales, which are the key to understanding the behavior of the algorithm.







## 7. Matrix Martingales

“Spielbank Wiesbaden,” Wikimedia Commons

Some of the text of this lecture is copied from my paper [Tro11a]. The treatment of corrector processes has not appeared before.

We plan to analyze a randomized, sequential algorithm that operates on matrices. For this purpose, we need to extend the theory of matrix concentration from independent sums to martingales. The purpose of this lecture is to present the main elements of this extension. In Lecture 8, we will require the full power of this approach.

### 7.1 Matrix-valued random processes

We begin with some basic definitions from the theory of random processes and their matrix-valued cousins.

#### 7.1.1 Martingales

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space, and let  $\mathcal{F}_0 \subset \mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots \subset \mathcal{F}$  be a filtration of the master sigma algebra. We write  $\mathbb{E}_k$  for the expectation conditioned on  $\mathcal{F}_k$ . That is,  $\mathbb{E}_k$  averages over all random choices that take place after the instant  $k$ .

A *martingale* is a (real-valued) random process  $\{Y_k : k = 0, 1, 2, \dots\}$  that is adapted to the filtration and that satisfies two properties:

$$\mathbb{E}_{k-1} Y_k = Y_{k-1} \quad \text{and} \quad \mathbb{E} |Y_k| < +\infty \quad \text{for } k = 1, 2, 3, \dots$$

For simplicity, we require the initial value of a martingale to be null:  $Y_0 = 0$ . The *difference sequence* is the random process defined by

$$X_k = Y_k - Y_{k-1} \quad \text{for } k = 1, 2, 3, \dots$$

Roughly, the present value of a martingale depends only on the past values, and the martingale has the *status quo* property: today, on average, is the same as yesterday.

A *supermartingale* is a (real-valued) random process  $\{S_k : k \geq 0\}$  that is adapted to the filtration and that satisfies

$$\mathbb{E}_{k-1} S_k \leq S_{k-1}.$$

In other words, a supermartingale is a process with diminishing expectations.

### 7.1.2 Matrix martingales

Matrix martingales are defined in much the same manner as scalar martingales. Consider a random process  $\{Y_k : k = 0, 1, 2, \dots\}$  whose values are matrices of finite dimension. We say that the process is a *matrix martingale* when  $Y_0 = \mathbf{0}$  and

$$\mathbb{E}_{k-1} Y_k = Y_{k-1} \quad \text{and} \quad \mathbb{E} \|Y_k\| < +\infty \quad \text{for } k = 1, 2, 3, \dots$$

We write  $\|\cdot\|$  for the *spectral norm* of a matrix, which returns its largest singular value. As before, we define the difference sequence  $\{X_k : k = 1, 2, 3, \dots\}$  via the relation

$$X_k = Y_k - Y_{k-1} \quad \text{for } k = 1, 2, 3, \dots$$

A matrix-valued random process is a martingale if and only if we obtain a scalar martingale when we track each fixed coordinate in time.

### 7.1.3 Adapted sequences

A sequence  $\{X_k\}$  of random matrices is *adapted* to the filtration when each  $X_k$  is measurable with respect to  $\mathcal{F}_k$ . That is,  $X_k$  is completely determined by random choices made up to and including instant  $k$ . We say that a sequence  $\{V_k\}$  of random matrices is *predictable* when each  $V_k$  is measurable with respect to  $\mathcal{F}_{k-1}$ . In particular, the sequence  $\{\mathbb{E}_{k-1} X_k\}$  of conditional expectations of an adapted sequence  $\{X_k\}$  is predictable. A *stopping time* is a random variable  $K : \Omega \rightarrow \mathbb{N}_0 \cup \{\infty\}$  that satisfies  $\{K \leq k\} \subset \mathcal{F}_k$  for  $k = 0, 1, 2, \dots, \infty$ .

### 7.1.4 Stopped processes

Suppose that  $\{S_k : k \geq 0\}$  is an adapted random process, and let  $K$  be a stopping time. The *stopped process*  $\{S_{k \wedge K} : k \geq 0\}$  coincides with the original process up to the stopping time  $K$ , after which it remains constant.

**Fact 7.1 (Stopped processes).** Let  $\{S_k\}$  be a (super)martingale, and let  $K$  be a stopping time. The stopped process  $\{S_{k \wedge K}\}$  remains a (super)martingale. ■

## 7.2 Tail bounds for matrix-valued processes

Now, let us develop a general methodology for establishing tail bounds for matrix-valued random processes. The basic technique can be traced at least as far as Freedman's work [Fre75] on scalar random processes. In the next section, we introduce the extra tools that are required to apply these results fruitfully in the matrix setting.

### 7.2.1 Corrector processes

We begin with the definition of a corrector process for a martingale. The corrector process is an auxiliary random process that provides an evolving bound on the growth of the martingale. This concept is rather abstract, but we will soon see how to make it more concrete.

**Definition 7.2 (Corrector process).** Let  $g : [0, \infty] \rightarrow [0, \infty]$  be a function. Consider a martingale  $\{Y_k : k = 0, 1, 2, \dots\}$  and a predictable random process  $\{W_k : k = 0, 1, 2, \dots\}$  that consist of self-adjoint random matrices with dimension  $d$ . Define the real-valued random processes

$$S_k(\theta) = \text{tr} \exp(\theta Y_k - g(\theta) W_k) \quad \text{for } \theta \geq 0. \quad (7.1)$$

We say that  $\{g W_k\}$  is a *corrector process* for the martingale  $\{Y_k\}$  if  $S_k(\theta)$  is a positive supermartingale for all  $\theta \geq 0$ .

Since we are assuming that the martingale has a null initial value ( $Y_0 = \mathbf{0}$ ), it is natural to require that the corrector process also has null initial value ( $W_0 = \mathbf{0}$ ). In this case, the initial value of the supermartingale satisfies  $S_0(\theta) = d$  for all  $\theta \geq 0$ . Furthermore, the supermartingale only takes positive values.

### 7.2.2 Lower bounds for the supermartingale

Next, we present a simple inequality that bounds the supermartingale  $S_k$  below when we have control on the eigenvalues of the two processes.

**Lemma 7.3** Suppose that  $\lambda_{\max}(Y) \geq t$  and that  $\lambda_{\max}(W) \leq w$ . For each  $\theta > 0$ ,

$$\text{tr} \exp(\theta Y - g(\theta) W) \geq e^{\theta t - g(\theta) w}.$$

*Proof.* Recall that  $g(\theta) \geq 0$ . The bound results from a calculation:

$$\begin{aligned} \text{tr} e^{\theta Y - g(\theta) W} &\geq \text{tr} e^{\theta Y - g(\theta) w \mathbf{I}} \\ &\geq \lambda_{\max} \left( e^{\theta Y - g(\theta) w \mathbf{I}} \right) = e^{\theta \lambda_{\max}(Y) - g(\theta) w} \geq e^{\theta t - g(\theta) w}. \end{aligned}$$

The first inequality depends on the semidefinite relation  $W \preceq w \mathbf{I}$  and the monotonicity of the trace exponential with respect to the semidefinite order (Fact 1.8). The second inequality relies on the fact that the trace of a psd matrix is at least as large as its maximum eigenvalue. The third identity follows from the spectral mapping theorem and elementary properties of the maximum eigenvalue map. ■

### 7.2.3 A tail bound for matrix martingales

Our key theorem provides a bound on the probability that the maximum eigenvalue of a matrix martingale ever exceeds a threshold.

**Theorem 7.4 (Master tail bound for matrix martingales).** Consider a matrix martingale  $\{Y_k\}$  consisting of self-adjoint matrices with dimension  $d$ . For a function  $g : [0, \infty] \rightarrow [0, \infty]$ , assume that  $\{g W_k\}$  is a corrector process for the martingale.

Then, for all  $t, w \in \mathbb{R}$ ,

$$\mathbb{P} \{ \exists k \geq 0 : \lambda_{\max}(\mathbf{Y}_k) \geq t \text{ and } \lambda_{\max}(\mathbf{W}_k) \leq w \} \leq d \cdot \inf_{\theta > 0} e^{-\theta t + g(\theta) w}.$$

*Proof.* The overall proof strategy is the same as the stopping-time technique used by Freedman [Fre75]. Fix a positive parameter  $\theta$ , which we will optimize later. Introduce the supermartingale  $S_k = S_k(\theta)$ , as in (7.1).

Define a stopping time  $K$  by finding the first time instant  $k$  when the maximum eigenvalue of the martingale reaches the level  $t$  even though the corrector process has maximum eigenvalue no larger than  $w$ . That is,

$$K := \inf \{ k \geq 0 : \lambda_{\max}(\mathbf{Y}_k) \geq t \text{ and } \lambda_{\max}(\mathbf{W}_k) \leq w \}.$$

When the infimum is empty, the stopping time  $K = \infty$ . Consider a system of exceptional events:

$$\mathbf{E}_k := \{ \lambda_{\max}(\mathbf{Y}_k) \geq t \text{ and } \lambda_{\max}(\mathbf{W}_k) \leq w \} \quad \text{for } k = 0, 1, 2, \dots$$

Construct the event  $\mathbf{E} := \bigcup_{k=0}^{\infty} \mathbf{E}_k$  that one or more of these exceptional situations takes place. The intuition behind this definition is that our control on the corrector process  $\{\mathbf{W}_k\}$  prevents the martingale  $\{\mathbf{Y}_k\}$  from exhibiting a large value. As a result, the event  $\mathbf{E}$  is rather unlikely.

We are prepared to estimate the probability of the exceptional event. First, note that  $K < \infty$  on the event  $\mathbf{E}$ . Therefore, Lemma 7.3 provides a conditional lower bound for the supermartingale  $\{S_k\}$  at the stopping time  $K$ :

$$S_K = \text{tr} \exp(\theta \mathbf{Y}_K - g(\theta) \mathbf{W}_K) \geq e^{\theta t - g(\theta) w} \quad \text{on the event } \mathbf{E}.$$

The stopped process  $\{S_{k \wedge K}\}$  is also a positive supermartingale with initial value  $d$ , so

$$d \geq \liminf_{k \rightarrow \infty} \mathbb{E}[S_{k \wedge K}] \geq \liminf_{k \rightarrow \infty} \mathbb{E}[S_{k \wedge K} \mathbb{1}_{\mathbf{E}}] \geq \mathbb{E}[\liminf_{k \rightarrow \infty} S_{k \wedge K} \mathbb{1}_{\mathbf{E}}] = \mathbb{E}[S_K \mathbb{1}_{\mathbf{E}}].$$

The indicator function decreases the expectation because the stopped process is positive. Fatou's lemma justifies the third inequality, and we have identified the limit using the fact that  $K < \infty$  on the event  $\mathbf{E}$ . It follows that

$$d \geq \mathbb{E}[S_K \mathbb{1}_{\mathbf{E}}] \geq (\mathbb{P}\mathbf{E}) \cdot \inf_{\mathbf{E}} S_K \geq (\mathbb{P}\mathbf{E}) \cdot e^{\theta t - g(\theta) w}.$$

Rearrange the relation to obtain

$$\mathbb{P}\mathbf{E} \leq d \cdot e^{-\theta t + g(\theta) w}.$$

Minimize the right-hand side with respect to  $\theta$  to complete the main part of the argument. ■

### 7.3 Building a corrector process

To convert Theorem 7.4 into a useful tool, we need a mechanism for constructing a corrector process. Fortunately, as in the case of independent sums of random matrices, Lieb's theorem comes to our rescue. We will see that we can construct a corrector process using matrix cgfs.

### 7.3.1 Correctors

Let us specialize the notion of a corrector to a single matrix. This will be the building block for constructing a corrector process.

**Definition 7.5 (Corrector).** Let  $g : [0, \infty] \rightarrow [0, \infty]$  be a function. Consider a random self-adjoint matrix  $\mathbf{X}$  and a fixed matrix  $\mathbf{V}$ , each with dimension  $d$ . We say that  $g \mathbf{V}$  is a *corrector* for  $\mathbf{X}$  when

$$\mathbb{E} \operatorname{tr} \exp(\mathbf{M} + \theta \mathbf{X} - g(\theta) \mathbf{V}) \leq \operatorname{tr} \exp(\mathbf{M}) \quad \text{for } \theta > 0.$$

This bound must hold for every fixed matrix  $\mathbf{M} \in \mathbb{H}_d$ .

### 7.3.2 Lieb's theorem and Tropp's corollary

Our main tool for producing explicit correctors is Lieb's theorem [Lie73, Thm. 6]. We refer to [Tro15, Chap. 8] for a digestible proof of this result.

**Theorem 7.6 (Lieb, 1973).** Fix a self-adjoint matrix  $\mathbf{H}$ . The function

$$\mathbf{A} \mapsto \operatorname{tr} \exp(\mathbf{H} + \log \mathbf{A})$$

is concave on the pd cone.

Lieb's theorem tells us that we can construct a corrector from a cumulant generating function. This simple but powerful observation first appeared in [Tro11a].

**Corollary 7.7 (Tropp, 2010).** Let  $\mathbf{M}$  be a fixed self-adjoint matrix, and let  $\mathbf{X}$  be a random self-adjoint matrix of the same dimension. For any  $\theta \in \mathbb{R}$ ,

$$\mathbb{E} \operatorname{tr} \exp(\mathbf{M} + \theta \mathbf{X} - \log \mathbb{E} e^{\theta \mathbf{X}}) \leq \operatorname{tr} \exp(\mathbf{M}).$$

*Proof.* Define the random matrix  $\mathbf{Y} = e^{\theta \mathbf{X}}$ , and calculate that

$$\begin{aligned} \mathbb{E} \operatorname{tr} \exp(\mathbf{M} + \theta \mathbf{X} - \log \mathbb{E} e^{\theta \mathbf{X}}) &= \mathbb{E} \operatorname{tr} \exp(\mathbf{M} + \log(\mathbf{Y}) - \log(\mathbb{E} \mathbf{Y})) \\ &\leq \operatorname{tr} \exp(\mathbf{M} + \log(\mathbb{E} \mathbf{Y}) - \log(\mathbb{E} \mathbf{Y})) \\ &= \operatorname{tr} \exp(\mathbf{M}). \end{aligned}$$

The first identity follows because the logarithm of the pd matrix  $\mathbf{Y}$  can be defined as the functional inverse of the matrix exponential. Theorem 7.6, with the fixed matrix  $\mathbf{H} = \mathbf{M} - \log(\mathbb{E} \mathbf{Y})$ , establishes that the trace function is concave in  $\mathbf{Y}$ . Invoke Jensen's inequality to draw the expectation inside the logarithm. ■

### 7.3.3 Example: The Bernstein corrector

Corollary 7.7 and Lemma 1.10 allow us to derive a corrector for a bounded, centered random matrix.

**Proposition 7.8 (Bernstein corrector).** Let  $\mathbf{X}$  be a random matrix that satisfies  $\mathbb{E} \mathbf{X} = \mathbf{0}$  and  $\|\mathbf{X}\| \leq 1$ . Then the matrix  $g(\theta)(\mathbb{E} \mathbf{X}^2)$  is a corrector for  $\mathbf{X}$ , where  $g(\theta) = (\theta^2/2)/(1 - |\theta|/3)$ .

*Proof.* We may calculate that

$$\mathbb{E} \operatorname{tr} \exp (\mathbf{M} + \theta \mathbf{X} - g(\theta)(\mathbb{E} \mathbf{X}^2)) \leq \mathbb{E} \operatorname{tr} \exp (\mathbf{M} + \theta \mathbf{X} - \log \mathbb{E} e^{\theta \mathbf{X}}) \leq \operatorname{tr} \exp (\mathbf{M}).$$

The first inequality follows from the Bernstein cgf bound, Lemma 1.10, because the trace exponential is monotone with respect to the semidefinite order (Fact 1.8). The second inequality is Corollary 7.7. ■

### 7.3.4 Example: The Chernoff corrector

Corollary 7.7 and Lemma 1.12 allow us to derive a corrector for a bounded, psd random matrix.

**Proposition 7.9 (Chernoff corrector).** Let  $\mathbf{X}$  be a random matrix that satisfies the bounds  $\mathbf{0} \preceq \mathbf{X} \preceq \mathbf{I}$ . Then the matrix  $g(\theta)(\mathbb{E} \mathbf{X})$  is a corrector for  $\mathbf{X}$ , where  $g(\theta) = e^\theta - 1$ .

We omit the repetitive proof.

### 7.3.5 From correctors to corrector processes

There is a straightforward connection between the corrector of a single random and the corrector process of a martingale.

**Proposition 7.10 (Corrector processes).** Fix a function  $g : [0, \infty] \rightarrow [0, \infty]$ . Let  $\{\mathbf{Y}_k\}$  be a self-adjoint matrix martingale with difference sequence  $\{\mathbf{X}_k\}$ . Let  $\{\mathbf{V}_k\}$  be a predictable sequence of self-adjoint matrices. For each  $k$ , suppose that  $g \mathbf{V}_k$  is a corrector for  $\mathbf{X}_k$ , conditional on  $\mathcal{F}_{k-1}$ . Then the predictable process

$$\mathbf{W}_k = \sum_{i=1}^k \mathbf{V}_i$$

generates a corrector  $\{g \mathbf{W}_k\}$  for the martingale  $\{\mathbf{Y}_k\}$ .

*Proof.* As above, define

$$S_k(\theta) = \operatorname{tr} \exp(\theta \mathbf{Y}_k - g(\theta) \mathbf{W}_k).$$

To prove that the process is a supermartingale, we follow a short chain of inequalities. Split off the last term from  $\mathbf{Y}_k$  and  $\mathbf{W}_k$  to see that

$$\begin{aligned} \mathbb{E}_{k-1} S_k(\theta) &= \mathbb{E}_{k-1} \operatorname{tr} \exp(\theta \mathbf{Y}_{k-1} - g(\theta) \mathbf{W}_{k-1} + \theta \mathbf{X}_k - g(\theta) \mathbf{V}_k) \\ &\leq \operatorname{tr} \exp(\theta \mathbf{Y}_{k-1} - g(\theta) \mathbf{W}_{k-1}) = S_{k-1}(\theta). \end{aligned}$$

This inequality follows immediately from the assumption that  $g \mathbf{V}_k$  is a corrector for  $\mathbf{X}_k$ , conditional on  $\mathcal{F}_{k-1}$ . We can apply this hypothesis because  $\mathbf{Y}_{k-1}$  and  $\mathbf{W}_{k-1}$  both are measurable with respect to  $\mathcal{F}_{k-1}$ . ■

### 7.3.6 Correctors tensorize

Let us continue with some general methods for constructing correctors of more complicated matrices. First, correctors tensorize over independent random matrices.



**Proposition 7.11 (Correctors tensorize).** Let  $g : [0, \infty] \rightarrow [0, \infty]$ . Consider an independent family  $\{X_k : k = 1, \dots, n\}$  of self-adjoint random matrices, and a nonrandom family  $\{V_k : k = 1, \dots, n\}$  of self-adjoint matrices. Suppose that  $g V_k$  is a corrector for  $X_k$  for each  $k$ . Then  $g \sum_k V_k$  is a corrector for  $\sum_k X_k$ .

*Proof.* This result follows by iteration of Definition 7.5. Let  $M$  be a fixed matrix.

$$\begin{aligned} \mathbb{E} \mathbb{E}_n \operatorname{tr} \exp \left( M + \theta \sum_{k=1}^n X_k - g(\theta) \sum_{k=1}^n V_k \right) \\ \leq \mathbb{E} \mathbb{E}_{n-1} \operatorname{tr} \exp \left( M + \theta \sum_{k=1}^{n-1} X_k - g(\theta) \sum_{k=1}^{n-1} V_k \right) \\ \leq \dots \leq \operatorname{tr} \exp(M). \end{aligned}$$

This is what we needed to show. ■

### 7.3.7 The composition rule

Next, let us present a composition rule that allows us to derive a corrector for a random matrix that is constructed in multiple steps.

**Proposition 7.12 (Composition rule).** Consider sigma fields  $\mathcal{F}_0 \subset \mathcal{F}_1 \subset \mathcal{F}_2$ . Let  $X$  be a random matrix that is measurable with respect to  $\mathcal{F}_2$ . For  $\theta \geq 0$ , suppose that

$$\begin{aligned} \mathbb{E} \left[ \operatorname{tr} \exp (M_1 + \theta X - g(\theta) V_1) \mid \mathcal{F}_1 \right] &\leq \operatorname{tr} \exp (M_1) \\ \mathbb{E} \operatorname{tr} \exp (M_0 + \theta V_1 - h(\theta) V_0) &\leq \operatorname{tr} \exp (M_0). \end{aligned}$$

In this expression,  $V_1$  and  $M_1$  are measurable with respect to  $\mathcal{F}_1$ , while  $V_0$  and  $M_0$  are measurable with respect to  $\mathcal{F}_0$ . Then  $(h \circ g) V_0$  is a corrector for  $X$ .

*Proof.* Let  $M$  be measurable with respect to  $\mathcal{F}_0$ . Calculate that

$$\begin{aligned} \mathbb{E} \operatorname{tr} \exp (M + \theta X - (h \circ g)(\theta) V_0) \\ = \mathbb{E} \mathbb{E} \left[ \operatorname{tr} \exp (M + \theta X - g(\theta) V_1 + g(\theta) V_1 - (h \circ g)(\theta) V_0) \mid \mathcal{F}_1 \right] \\ \leq \mathbb{E} \operatorname{tr} \exp (M + g(\theta) V_1 - h(g(\theta)) V_0) \\ \leq \operatorname{tr} \exp(M). \end{aligned}$$

This is the definition of a corrector. ■

## 7.4 Example: The matrix Freedman inequality

As an example, let us use Theorem 7.4 to prove the matrix version of a classic martingale inequality due to Freedman.

**Theorem 7.13 (Matrix Freedman).** Consider a matrix martingale  $\{Y_k\}$  consisting of self-adjoint matrices with dimension  $d$ . Assume that the difference sequence  $\{X_k\}$  satisfies

$$\|X_k\| \leq B \quad \text{for } k = 1, 2, 3, \dots$$



Define the cumulative predictable quadratic variation process:

$$\mathbf{W}_0 = \mathbf{0} \quad \text{and} \quad \mathbf{W}_k = \sum_{i=1}^k \mathbb{E}_{i-1} \mathbf{X}_i^2 \quad \text{for } k = 1, 2, 3, \dots$$

Then, for all  $t \geq 0$  and  $\sigma^2 \geq 0$ ,

$$\mathbb{P} \left\{ \exists k \geq 0 : \lambda_{\max}(\mathbf{Y}_k) \geq t \quad \text{and} \quad \lambda_{\max}(\mathbf{W}_k) \leq \sigma^2 \right\} \leq d \cdot \exp \left( \frac{-t^2/2}{\sigma^2 + Bt/3} \right).$$

*Proof.* We assume that  $B = 1$ ; the general result follows by re-scaling since  $\mathbf{Y}_k$  is 1-homogeneous and  $\mathbf{W}_k$  is 2-homogeneous.

Invoke Proposition 7.8 conditionally to obtain a corrector process for  $\{\mathbf{X}_k\}$ . Indeed, we can choose

$$\mathbf{V}_k = \mathbb{E}_{k-1} \mathbf{X}_k^2 \quad \text{and} \quad g(\theta) = \frac{\theta^2/2}{1 - |\theta|/3}.$$

Theorem 7.4 now implies that

$$\mathbb{P} \left\{ \exists k \geq 0 : \lambda_{\max}(\mathbf{Y}_k) \geq t \quad \text{and} \quad \lambda_{\max}(\mathbf{W}_k) \leq \sigma^2 \right\} \leq d \cdot \inf_{\theta > 0} e^{-\theta t + g(\theta) \sigma^2}.$$

Make the inspired choice  $\theta = t/(\sigma^2 + t/3)$  to complete the argument. ■

**Exercise 7.1** Extend the matrix Freedman inequality to a martingale sequence consisting of rectangular matrices.



## 8. Sparse Cholesky

©COMSOL Multiphysics

This lecture is adapted from Rasmus Kyng's thesis [Kyn17]. The application of matrix martingales has been streamlined by using the notion of a corrector process.

In this lecture, we introduce a practical algorithm for solving Laplacian linear systems in near-linear time. The algorithm is remarkable in its simplicity, but the analysis relies on many of the sophisticated ideas that we have encountered in the previous lectures.

This approach, called the SparseCholesky algorithm, was developed by Rasmus Kyng and Sushant Sachdeva [KS16]. It was further refined in Kyng's dissertation [Kyn17]. It is closely related to an earlier algorithm for connection Laplacians, developed by Dan Spielman's group [Kyn+16]. Altogether, these methods hold real promise for solving large graph Laplacian systems in practice.

### 8.1 Approximate solutions of Laplacian systems

We begin with a high-level approach for computing an approximation solution of a Laplacian system via preconditioning.

#### 8.1.1 Approximate solutions

Let  $L$  be the Laplacian matrix of a connected multigraph. Suppose that we wish to find the unique solution  $x_\star$  to the linear system

$$Lx = f \quad \text{where} \quad \mathbf{1}^* f = 0 \quad \text{and} \quad \mathbf{1}^* x = 0.$$

For a parameter  $\varepsilon > 0$ , we can relax our requirement by asking for an approximate solution  $\mathbf{x}_\varepsilon$  that satisfies the relative error bound

$$\|\mathbf{x}_\varepsilon - \mathbf{x}_\star\|_L \leq \varepsilon \cdot \|\mathbf{x}_\star\|_L$$

Here,  $\|\cdot\|_L$  is the norm associated with the quadratic form (i.e., the Dirichlet form) determined by  $\mathbf{L}$ . That is,  $\|\mathbf{x}\|_L = (\mathbf{x}^* \mathbf{L} \mathbf{x})^{1/2}$ .

### 8.1.2 Approximate Cholesky decomposition

Suppose that we are able to construct a sparse, approximate Cholesky decomposition of the Laplacian matrix:

$$0.5 \mathbf{L} \preceq \mathbf{C} \mathbf{C}^* \preceq 1.5 \mathbf{L} \quad \text{where} \quad \text{nnz}(\mathbf{C}) = O(m \log n). \quad (8.1)$$

The symbol  $\preceq$  refers to the semidefinite order. The matrix  $\mathbf{C}$  is morally lower-triangular; in other words, there is a permutation of coordinates that brings the matrix into lower-triangular form. The function  $\text{nnz}$  returns the number of nonzero entries in a matrix.

### 8.1.3 Preconditioning

Given the sparse, approximate factor  $\mathbf{C}$ , we can precondition the linear system (6.1):

$$(\mathbf{C}^\dagger \mathbf{L} \mathbf{C}^{*\dagger})(\mathbf{C}^* \mathbf{x}) = (\mathbf{C}^\dagger \mathbf{f}).$$

Owing to (8.1), the preconditioned linear system has condition number  $\kappa \leq 3$ . Of course, in practice, we treat the matrix as an operator acting on vectors. Each time we apply the operator, we use forward and back substitution to invoke  $\mathbf{C}^\dagger$  and  $\mathbf{C}^{*\dagger}$ . The total cost of each application of the matrix is  $\Theta(m \log n)$  arithmetic operations, because the substitution method for solving a morally triangular system exploits sparsity.

We can solve the preconditioned system using the conjugate gradient algorithm. If the initial iterate  $\mathbf{x}_0 = \mathbf{0}$ , then, after  $j$  iterations, we attain the error bound

$$\|\mathbf{x}_j - \mathbf{x}_\star\|_L \leq 2 \left[ \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right]^j \|\mathbf{x}_\star\|_L \leq 3^{1-j} \|\mathbf{x}_\star\|_L.$$

In particular, we can achieve a relative error of  $\varepsilon$  in the Dirichlet energy norm after  $O(\log(1/\varepsilon))$  iterations.

### 8.1.4 Summary

In summary, once we have constructed an approximate Cholesky decomposition that satisfies (8.1), we can solve the linear system to relative error  $\varepsilon$  using  $O(m \log(n) \log(1/\varepsilon))$  arithmetic operations. This computation takes place in time that is nearly linear in the number of degrees of freedom in the graph.

Easy! We just have to achieve (8.1). In the rest of this lecture, we will explain how to perform this feat.

## 8.2 Overview of the algorithm

Let us begin with an overview of the **SparseCholesky** algorithm for computing a sparse, approximate Cholesky decomposition that satisfies (8.1). We will fill in the details and perform the analysis over the balance of the lecture.

### 8.2.1 Setup

Fix the ground set  $V = \{1, \dots, n\}$  of vertices. Let  $G$  be a connected multigraph on  $V$  composed of  $m$  weighted multiedges. As usual, we will interact with the multigraph  $G$  via its Laplacian matrix  $L$ . The Laplacian will be represented as a sum over multiedges:

$$L = \sum_{e \in L} w_L(e) \Delta_e \in \mathbb{H}_V.$$

The indexing for the sum and the notation for the weight function are intended to be mnemonic, if not overly formal.

### 8.2.2 The SparseCholesky procedure

The **SparseCholesky** algorithm is based on the same template as the ordinary Cholesky decomposition, but it judiciously injects randomness to minimize the computational burden. The basic idea is to randomly sample the cliques that arise as we eliminate vertices from the multigraph:

#### Summary of **SparseCholesky**:

For each iteration  $i = 1, \dots, n$ :

1. Select a random vertex  $u_i$  to eliminate.
2. Add a random approximation of the clique induced by  $u_i$ .
3. Remove the star induced by  $u_i$ .

We continue with a more detailed presentation of the steps in the procedure.

#### Preprocessing

Before we begin, we split each multiedge into a fixed number of pieces to reduce the leverage of each multiedge below a threshold.

As the algorithm constructs new multiedges, we will ensure that the leverages never increase beyond the initial threshold. This property helps control the variance of the random clique approximations.

#### Initialization

Let  $S_0 = L$ . Let  $F_0 = V$  of vertices that have not been eliminated.

We will maintain the invariant that the iterate  $S_i$  is supported on the vertices listed in  $F_i$ . The number of vertices remaining at each step will satisfy  $|F_i| = n - i$ .

#### Selecting a vertex to eliminate

At each iteration  $i = 1, 2, 3, \dots, n$ , select a vertex  $u_i$  uniformly at random from  $F_{i-1}$ . Update  $F_i = F_{i-1} \setminus \{u_i\}$ .

Selecting a random vertex  $u_i$  renders it unlikely that there are many multiedges incident on  $u_i$ . Furthermore, it is unlikely that the clique induced by  $u_i$  involves multiedges whose total leverage is large. These facts are critical for controlling the runtime of the algorithm and ensuring that it produces an accurate approximation.

### Collecting information

Extract data from the current iterate  $\mathbf{S}_{i-1}$ :

$$\mathbf{c}_i = \frac{1}{\sqrt{(\mathbf{S}_{i-1})_{u_i u_i}}} \mathbf{S}_{i-1} \delta_{u_i}.$$

If  $(\mathbf{S}_{i-1})_{u_i u_i} = 0$ , then we set  $\mathbf{c}_i = \mathbf{0}$ . Since the matrix  $\mathbf{S}_{i-1}$  is supported on the coordinates listed in  $F_{i-1}$ , the support of  $\mathbf{c}_i$  is also contained in  $F_{i-1}$ .

### Sampling the clique

To proceed, we will approximate the Schur complement  $\mathbf{S}_{i-1}/u_i$ . To do so, we first construct a random sparse Laplacian matrix  $\mathbf{K}_i$  that approximates the clique induced by eliminating  $u_i$ . We will explain how to perform this approximation later, in Section 8.4. The basic requirement on that  $\mathbf{K}_i$  is that

$$\mathbb{E}_{i-1}[\mathbf{K}_i \mid u_i] = \text{CLIQUE}(u_i, \mathbf{S}_{i-1}).$$

The expectation  $\mathbb{E}_{i-1}$  conditions on all of the randomness in the first  $i - 1$  iterations. We also condition separately on the random vertex  $u_i$  drawn at step  $i$ .

The number of multiedges in the clique approximation  $\mathbf{K}_i$  will not exceed the total number of multiedges incident on the vertex  $u_i$ , so the number of multiedges remaining in the multigraph does not increase as the iteration advances. This property also ensures that the cost of computing the clique approximation is under control.

Moreover, we will ensure that the clique approximation  $\mathbf{K}_i$  has no multiedges incident on  $u_1, \dots, u_i$ . That is,  $\mathbf{K}_i$  is supported on the coordinates listed in  $F_i$ .

### Approximating the Schur complement

Now, form the approximate Schur complement:

$$\mathbf{S}_i = (\mathbf{S}_{i-1} - \text{STAR}(u_i, \mathbf{S}_{i-1})) + \mathbf{K}_i. \quad (8.2)$$

In the last step, we set  $\mathbf{S}_n = \mathbf{0}$ .

For reference, the star (6.2) and clique (6.3) induced by a vertex were defined before. This construction ensures that  $\mathbf{S}_i$  is supported on the coordinates listed in  $F_i$ . Therefore, we continue to reduce the size of the problem.

### Forming the decomposition

As usual, we conclude by compiling the matrix

$$\mathbf{C} = [\mathbf{c}_1 \quad \dots \quad \mathbf{c}_n] \in \mathbb{R}^{V \times V}.$$

By construction of the vectors  $\mathbf{c}_i$ , the matrix  $\mathbf{C}$  is morally lower-triangular. The elimination order is associated with the permutation  $\pi$  defined by  $\pi(i) = u_i$  for  $i = 1, \dots, n$ .

### 8.2.3 Laplacian approximations

How do we make sense of this approach? Note that the SparseCholesky iteration induces a sequence of approximations to the Laplacian matrix:

$$\mathbf{L}_i = \mathbf{S}_i + \sum_{k=1}^i \mathbf{c}_k \mathbf{c}_k^* \quad \text{for } i = 0, 1, 2, \dots, n.$$

In particular, the initial value of the sequence is the original Laplacian, while the final value is our approximate Cholesky decomposition:

$$\mathbf{L}_0 = \mathbf{L} \quad \text{and} \quad \mathbf{L}_n = \sum_{i=1}^n \mathbf{c}_i \mathbf{c}_i^* = \mathbf{C} \mathbf{C}^*.$$

The difference sequence of the random process  $\{\mathbf{L}_i\}$  satisfies

$$\begin{aligned} \mathbf{L}_i - \mathbf{L}_{i-1} &= \mathbf{S}_i - \mathbf{S}_{i-1} + \mathbf{c}_i \mathbf{c}_i^* \\ &= \mathbf{K}_i - \text{STAR}(u_i, \mathbf{S}_{i-1}) + \mathbf{c}_i \mathbf{c}_i^* \\ &= \mathbf{K}_i - \text{CLIQUE}(u_i, \mathbf{S}_{i-1}) \\ &= \mathbf{K}_i - \mathbb{E}[\mathbf{K}_i \mid u_i]. \end{aligned}$$

The second relation follows from the definition (8.2) of the approximate Schur complement  $\mathbf{S}_i$ . We have used an identity from the last lecture:

$$\text{CLIQUE}(u_i, \mathbf{S}_{i-1}) = \text{STAR}(u_i, \mathbf{S}_{i-1}) - \mathbf{c}_i \mathbf{c}_i^*. \quad (8.3)$$

It is now evident that each increment is conditionally zero mean:

$$\mathbb{E}_{i-1}[\mathbf{L}_i - \mathbf{L}_{i-1}] = \mathbb{E}_{i-1}[\mathbb{E}_{i-1}[\mathbf{L}_i - \mathbf{L}_{i-1} \mid u_i]] = \mathbf{0}.$$

In particular,

$$\mathbb{E} \mathbf{L}_i = \mathbf{L} \quad \text{for each } i = 1, 2, 3, \dots, n.$$

We discover that  $\{\mathbf{L}_i - \mathbf{L}_0\}$  is a matrix martingale with null initial value. The final value of this martingale is the error in the approximate Cholesky decomposition:

$$\mathbf{C} \mathbf{C}^* - \mathbf{L} = \mathbf{L}_n - \mathbf{L}_0 = \sum_{i=1}^n (\mathbf{L}_i - \mathbf{L}_{i-1}).$$

Therefore, we can use the theory of matrix martingales to understand the behavior of the algorithm.

## 8.3 Preliminaries for the analysis

Let us begin the argument with some preliminary notation and simplifications.

### 8.3.1 The normalizing map

We define the normalizing map  $\Phi$  associated with the Laplacian  $\mathbf{L}$  of the initial multigraph  $\mathbf{G}$ :

$$\Phi(\mathbf{M}) = (\mathbf{L}^\dagger)^{1/2} \mathbf{M} (\mathbf{L}^\dagger)^{1/2} \quad \text{for } \mathbf{M} \in \mathbb{H}_V.$$

This map has two properties that will play a role in the argument. Since  $\mathbf{G}$  is connected,  $\Phi(\mathbf{L}) = \mathbf{P}$ , where  $\mathbf{P}$  is the orthogonal projector onto  $\text{lin}\{\mathbf{1}\}^\perp$ . Second,  $\Phi$  is a *positive map*. That is,

$$\mathbf{M} \succeq \mathbf{0} \quad \text{implies} \quad \Phi(\mathbf{M}) \succeq \mathbf{0}.$$

Let us emphasize that  $\Phi$  is always constructed from the Laplacian  $\mathbf{L}$  of the initial multigraph.

### 8.3.2 The approximation requirement

Recall that the random process  $\{\mathbf{L}_i\}$  has the terminal value  $\mathbf{L}_n = \mathbf{C}\mathbf{C}^*$ . We can express the approximation requirement (8.1) as

$$-0.5 \mathbf{L} \preceq \mathbf{L}_n - \mathbf{L} \preceq +0.5 \mathbf{L}.$$

Since  $\mathbf{L}_n$  is a Laplacian, its range is contained in the range of the Laplacian  $\mathbf{L}$ . Therefore, we can apply the normalizing map to obtain an equivalent condition:

$$-0.5 \mathbf{P} \preceq \Phi(\mathbf{L}_n - \mathbf{L}) \preceq +0.5 \mathbf{P}.$$

Using the relation  $\mathbf{L}_0 = \mathbf{L}$  and taking care with the ranges of the matrices that appear, we can write the latter expression as a pair of eigenvalue bounds:

$$\begin{aligned} \lambda_{\max}(\Phi(\mathbf{L}_n - \mathbf{L}_0)) &\leq +0.5; \\ \lambda_{\min}(\Phi(\mathbf{L}_n - \mathbf{L}_0)) &\geq -0.5. \end{aligned}$$

In other words, we must control the discrepancy between the terminal value  $\mathbf{L}_n$  and the initial value  $\mathbf{L}_0$  of the random process. Matrix martingale inequalities are tailor-made for this purpose.

To see how this will work, we decompose the martingale into its difference sequence:

$$\begin{aligned} \Phi(\mathbf{L}_n - \mathbf{L}_0) &= \sum_{i=0}^n \Phi(\mathbf{L}_i - \mathbf{L}_{i-1}) \\ &= \sum_{i=0}^n \Phi(\mathbf{K}_i - \text{CLIQUE}(u_i, \mathbf{S}_{i-1})) \\ &= \sum_{i=0}^n \Phi(\mathbf{K}_i - \mathbb{E}_{i-1}[\mathbf{K}_i \mid u_i]). \end{aligned}$$

The next step is to construct and analyze the randomized clique estimators  $\mathbf{K}_i$ . As a result, we will obtain a corrector process for the martingale  $\{\mathbf{L}_i\}$ , which will lead to the required tail bounds.

### 8.3.3 Splitting the edges

Recall that many matrix concentration bounds, such as the matrix Bernstein and matrix Freedman inequalities, require some type of uniform control over the random contributions. To obtain this control, we will preprocess the multigraph by splitting each multiedge into pieces.

Let  $R \geq 1$  be a parameter that we will fix later. (To be concrete, we will set  $R = \Theta(\log n)$ .) We construct a new multigraph, with Laplacian  $\mathbf{L}'$ , by splitting each



edge in the Laplacian  $\mathbf{L}$  into  $R$  equal pieces. This action has the effect of multiplying each leverage by a factor of  $1/R$ .

More precisely, we iterate over each multiedge  $e = uv$  in  $\mathbf{L}$ ; its weight is denoted as  $w_L(e)$ . We augment the new Laplacian  $\mathbf{L}'$  with  $R$  edges:

$$e_j = uv \quad \text{with} \quad w_{\mathbf{L}'}(e_j) = \frac{1}{R} w_L(e) \quad \text{for each } j = 1, 2, 3, \dots, R.$$

As matrices, the Laplacians are equal:  $\mathbf{L}' = \mathbf{L}$ . Regarded as multigraphs,  $\mathbf{L}'$  now has  $M = Rm$  multiedges, whereas  $\mathbf{L}$  only has  $m$  multiedges.

The leverage of each multiedge with respect to the new multigraph satisfies

$$w_{\mathbf{L}'}(e_j) \varrho_{\mathbf{L}'}(u, v) = \frac{1}{R} w_L(e) \varrho_L(u, v) \leq \frac{1}{R}.$$

Indeed, since the Laplacians are equal, the effective resistance of each pair of vertices is the same in both graphs. The last identity holds because of Proposition 5.6. Every multiedge that we construct during the algorithm will satisfy this same bound.

To avoid an extra notational burden, we will simply assume that the input Laplacian  $\mathbf{L}$  consists of  $M = mR$  multiedges, each with leverage score bounded by  $1/R$ . Effective resistances  $\varrho$  will always be computed with respect to this Laplacian  $\mathbf{L}$ .

## 8.4 Sampling from a clique

The main challenge in the **SparseCholesky** algorithm is to avoid the cost of constructing the full clique when eliminating a vertex. As noted, we plan to accomplish this goal using randomized sampling. This section explains how to perform this task.

### 8.4.1 Setup

Let  $\mathbf{S}$  be the Laplacian of a weighted multigraph over the set  $\mathbf{V}$  of vertices. Let  $\mathbf{F}$  be the *support* of  $\mathbf{S}$ ; that is, the subset of vertices where  $\mathbf{S}$  has an incident edge.

We will make two strong assumptions. First, we will assume that each multiedge in  $\mathbf{S}$  has bounded leverage with respect to the original Laplacian:

$$\text{For } e = uv \in \mathbf{S}, \quad w_{\mathbf{S}}(e) \varrho(u, v) \leq \frac{1}{R}. \quad (8.4)$$

We will often subscript the Laplacian  $\mathbf{S}$  to specify the multigraph. Second, we will assume that

$$\|\Phi(\mathbf{S})\| \leq 2. \quad (8.5)$$

In other words, the entire multigraph specified by  $\mathbf{S}$  has bounded leverage with respect to the target Laplacian  $\mathbf{L}$ .

The Laplacian  $\mathbf{S}$  evolves as the **SparseCholesky** algorithm progresses. We will ensure that these two properties hold, by force if necessary.

### 8.4.2 Eliminating a vertex

Fix a vertex  $u$  to eliminate from the Laplacian  $\mathbf{S}$ . To do so, we first construct the star induced by the vertex:

$$\text{STAR}(u, \mathbf{S}) = \sum_{e=uv \in \mathbf{S}} w_{\mathbf{S}}(e) \Delta_e.$$

That is, the star contains all of the multiedges in  $\mathbf{S}$  that are incident on  $u$ . Recall that  $\deg(u, \mathbf{S})$  is the number of multiedges incident on  $u$  in  $\mathbf{S}$ , i.e., the cardinality of the star. The total weight of the vertex is

$$w_{\mathbf{S}}(u) = \sum_{e \in \text{STAR}(u, \mathbf{S})} w_{\mathbf{S}}(e).$$

The clique induced by  $u$  has the Laplacian matrix

$$\text{CLIQUE}(u, \mathbf{S}) = \sum_{\substack{e_1=uv_1 \in \text{STAR}(u, \mathbf{S}) \\ e_2=uv_2 \in \text{STAR}(u, \mathbf{S})}} \frac{w_{\mathbf{S}}(e_1) w_{\mathbf{S}}(e_2)}{2w_{\mathbf{S}}(u)} \Delta_{v_1 v_2}.$$

Recall that each multiedge  $e$  in the star appears once in each sum, so the total number of multiedges in the clique is  $\deg_{\mathbf{S}}(u)^2$ .

Our project is to construct a Laplacian matrix  $\mathbf{K}$  that serves as an unbiased estimator for the clique:

$$\mathbb{E}[\mathbf{K} \mid u] = \text{CLIQUE}(u, \mathbf{S}).$$

We will insist that each multiedge in the approximation  $\mathbf{K}$  has the form  $v_1 v_2$  where the multiedges  $e_1 = uv_1$  and  $e_2 = uv_2$  both appear in  $\text{STAR}(u, \mathbf{S})$ . Moreover, the total number of multiedges in  $\mathbf{K}$  will not exceed  $\deg_{\mathbf{S}}(u)$ , the number of multiedges in the star that we remove. This is a quadratic reduction in complexity!

### 8.4.3 The sampling procedure

We are now prepared to detail the method for constructing a sparse, unbiased estimator of the clique.

#### Summary of CliqueSample:

1. Construct a probability mass on the multiedges in the star:

$$p(e) = \frac{w_{\mathbf{S}}(e)}{w_{\mathbf{S}}(u)} \quad \text{for each } e \in \text{STAR}(u, \mathbf{S}).$$

2. Draw a random multiedge  $e_1 = uv_1$  from the multiedges in  $\text{STAR}(u, \mathbf{S})$  according to the probability mass  $\mathbf{p}$ .
3. Draw a second random multiedge  $e_2 = uv_2$  from the multiedges in  $\text{STAR}(u, \mathbf{S})$  according to the uniform distribution.

4. Form the random Laplacian matrix of a new multiedge:

$$\mathbf{X} = \frac{w_S(e_1) w_S(e_2)}{w_S(e_1) + w_S(e_2)} \Delta_{v_1 v_2}. \quad (8.6)$$

This construction has several important features that we will explore in the next paragraphs.

For now, we remark that this sampling procedure is analogous to the other matrix sampling approximations that we have discussed throughout the course. The closest parallel, naturally, is with the sparse graph approximation in Lecture 5. In that context, we sampled edges in proportion to their leverages. In the present context, we do not know the leverage scores. Instead, we exploit the fact that the effective resistances satisfy a triangle inequality to obtain adequate sampling probabilities.

#### 8.4.4 Expectation of the random multiedge

First, let us compute the expectation of the Laplacian  $\mathbf{X}$  of a random multiedge. We will see that the random multiedge is an unbiased estimator of the clique, up to a fixed scale factor.

**Proposition 8.1 (Expectation of random Laplacian).** The Laplacian  $\mathbf{X}$  of the random multiedge (8.6) satisfies

$$\mathbb{E} \mathbf{X} = \frac{1}{\deg_S(u)} \cdot \text{CLIQUE}(u, \mathbf{S}).$$

*Proof.* This result follows by direct calculation. Below, each of the sums iterates over the multiedges in  $\text{STAR}(u, \mathbf{S})$ , which we omit from the notation.

$$\begin{aligned} \mathbb{E} \mathbf{X} &= \sum_{e_1=uv_1} \frac{w_S(e_1)}{w_S(u)} \sum_{e_2=uv_2} \frac{1}{\deg_S(u, \mathbf{S})} \cdot \frac{w_S(e_1) w_S(e_2)}{w_S(e_1) + w_S(e_2)} \Delta_{v_1 v_2} \\ &= \frac{1}{\deg_S(u)} \sum_{\substack{e_1=uv_1 \\ e_2=uv_2}} \frac{w_S(e_1) w_S(e_2)}{w_S(u)} \cdot \frac{w_S(e_1)}{w_S(e_1) + w_S(e_2)} \Delta_{v_1 v_2} \\ &= \frac{1}{\deg_S(u)} \sum_{\substack{e_1=uv_1 \\ e_2=uv_2}} \frac{w_S(e_1) w_S(e_2)}{2w_S(u)} \Delta_{v_1 v_2} = \frac{\text{CLIQUE}(u, \mathbf{S})}{\deg_S(u)}. \end{aligned}$$

The passage to the last line follows from the symmetry of the summands with respect to  $v_1$  and  $v_2$ . ■

#### 8.4.5 Each multiedge has bounded leverage

Next, let us verify that the multiedge  $\mathbf{X}$  constructed in (8.6) still has bounded leverage.

**Proposition 8.2 (Bounded leverage).** The random Laplacian matrix  $\mathbf{X}$  defined in (8.6) satisfies the uniform bound

$$\|\Phi(\mathbf{X})\| \leq \frac{1}{R}.$$

Equivalently, given multiedges  $e_1 = uv_1$  and  $e_2 = uv_2$ , the multiedge  $e = v_1v_2$  with weight

$$w_e = \frac{w_S(e_1) w_S(e_2)}{w_S(e_1) + w_S(e_2)}$$

has leverage score

$$w_e \varrho(v_1, v_2) \leq \frac{1}{R}.$$

*Proof.* This result is a consequence of the triangle inequality for effective resistances, Theorem 5.2. Indeed,

$$\begin{aligned} & w_S(e_1) w_S(e_2) \cdot \varrho(v_1, v_2) \\ & \leq w_S(e_2) \cdot w_S(e_1) \varrho(u, v_1) + w_S(e_1) \cdot w_S(e_2) \varrho(u, v_2) \\ & \leq \frac{1}{R} [w_S(e_1) + w_S(e_2)]. \end{aligned}$$

The last inequality holds because the weighted multiedges in  $S$  satisfy the uniform bound (8.4). Divide through by the bracket and identify the multiedge  $e = v_1v_2$  with weight  $w_e$  to arrive at the stated result. ■

#### 8.4.6 Corrector for the random multiedge

We are now prepared to bound the corrector for the Laplacian  $X$  of the random multiedge (8.6). First, we center and normalize the random matrix. The result is then an immediate application of the Bernstein corrector construction, Proposition 7.8.

**Proposition 8.3 (Corrector of random multiedge).** Fix a vertex  $u$ . The random matrix  $\Phi(X - \mathbb{E} X)$  admits the corrector

$$g(\theta) \cdot \frac{\Phi(\text{CLIQUE}(u, S))}{\deg_S(u)} \quad \text{where} \quad g(\theta) = \frac{\theta^2/(2R)}{1 - |\theta|/(3R)}.$$

The random matrix  $X$  is defined in (8.6). Let us emphasize that the vertex  $u$  is not random at this stage.

*Proof.* Proposition 8.1 implies that the random matrix  $\Phi(X - \mathbb{E} X)$  has mean zero. We have the uniform norm bound

$$\begin{aligned} \|\Phi(X - \mathbb{E} X)\| &= \|\Phi(X) - \Phi(\mathbb{E} X)\| \\ &= \max\{\|\Phi(X)\|, \|\Phi(\mathbb{E} X)\|\} \leq \|\Phi(X)\| \leq \frac{1}{R}. \end{aligned}$$

Since  $X$  is psd, so are  $\Phi(X)$  and  $\Phi(\mathbb{E} X)$ . This justifies the norm identity. The first inequality is Jensen's. The second inequality is Proposition 8.2.

Let us compute the variance:

$$\begin{aligned} \mathbb{E} \Phi(X - \mathbb{E} X)^2 &= \mathbb{E} \Phi(X)^2 - \Phi(\mathbb{E} X)^2 \leq \mathbb{E} \Phi(X)^2 \\ &\leq \mathbb{E} [\|\Phi(X)\| \cdot \Phi(X)] \leq \frac{1}{R} \Phi(\mathbb{E} X) \\ &= \frac{\Phi(\text{CLIQUE}(u, S))}{R \deg_S(u)}. \end{aligned}$$

We have repeatedly used the fact that  $\Phi$  is a positive linear map. The last identity is Proposition 8.1.

We now arrive at the result as a consequence of the Bernstein corrector bound, Proposition 7.8, and a scaling argument. ■

#### 8.4.7 An unbiased estimator for the clique

Next, we construct an estimator for the clique induced by eliminating the fixed vertex  $u$  from the Laplacian  $\mathbf{S}$ . To do so, add up  $\deg_{\mathbf{S}}(u)$  independent copies of the random Laplacian  $\mathbf{X}$  defined in (8.6):

$$\mathbf{K} = \sum_{j=1}^{\deg_{\mathbf{S}}(u)} \mathbf{X}_j \quad \text{where } \mathbf{X}_j \sim \mathbf{X} \text{ iid.} \quad (8.7)$$

Since Laplacians form a convex cone,  $\mathbf{K}$  is also the Laplacian of a multigraph. Let us verify that  $\mathbf{K}$  is an unbiased estimator of the clique and compute its corrector.

**Proposition 8.4 (Corrector of clique estimator: Fixed vertex).** Fix a vertex  $u$ . The random matrix  $\mathbf{K}$  defined in (8.7) is an unbiased estimator of the clique induced by eliminating  $u$  from  $\mathbf{S}$ :

$$\mathbb{E} \mathbf{K} = \text{CLIQUE}(u, \mathbf{S}).$$

The centered random matrix  $\Phi(\mathbf{K} - \mathbb{E} \mathbf{K})$  has corrector

$$g(\theta) \cdot \Phi(\text{CLIQUE}(u, \mathbf{S})) \quad \text{with} \quad g(\theta) = \frac{\theta^2/(2R)}{1 - |\theta|/(3R)}.$$

As before, we treat the vertex  $u$  as nonrandom.

*Proof.* The centered clique estimator decomposes as an independent sum:

$$\mathbf{K} - \mathbb{E} \mathbf{K} = \sum_{i=1}^{\deg_{\mathbf{S}}(u)} (\mathbf{X}_i - \mathbb{E} \mathbf{X}_i).$$

Proposition 7.11 states that the corrector tensorizes (over an independent sum). The result follows from Proposition 8.3. ■

#### 8.4.8 The clique induced by a random vertex

To complete our analysis of clique sampling, we consider what happens when we draw the vertex  $u$  at random.

First, let us develop several properties of the clique induced by an arbitrary vertex  $u$ . Recall from (8.3) that  $\text{CLIQUE}(u, \mathbf{S})$  is a Laplacian matrix obtained by subtracting a psd matrix from  $\text{STAR}(u, \mathbf{S})$ . Moreover,  $\text{STAR}(u, \mathbf{S})$  is the Laplacian of a subset of multiedges in  $\mathbf{S}$ . Therefore,

$$\mathbf{0} \preceq \text{CLIQUE}(u, \mathbf{S}) \preceq \text{STAR}(u, \mathbf{S}) \preceq \mathbf{S}.$$

Using the assumption (8.5), we obtain the bound

$$\|\Phi(\text{CLIQUE}(u, \mathbf{S}))\| \leq \|\Phi(\mathbf{S})\| \leq 2.$$

In other words, the whole clique has bounded leverage.

Second, we compute the average of the clique with respect to a vertex  $u$  drawn uniformly from the support  $F$  of the Laplacian  $S$ . Note that

$$\begin{aligned}\mathbb{E}_u \text{CLIQUE}(u, S) &\leq \mathbb{E}_u \text{STAR}(u, S) = \frac{1}{|F|} \cdot \sum_{u \in F} \sum_{e \in \text{STAR}(u, S)} w_S(e) \Delta_e \\ &= \frac{2}{|F|} \sum_{e \in S} w_S(e) \Delta_e = \frac{2}{|F|} \cdot S.\end{aligned}$$

Indeed, every multiedge in  $S$  appears twice in the sum because we touch each of its two endpoints as we loop over the vertices in the support  $F$  of  $S$ . Applying the normalizing map,

$$\mathbb{E}_u \Phi(\text{CLIQUE}(u, S)) \leq \frac{2}{|F|} \cdot \Phi(S) \leq \frac{4}{|F|} \cdot I.$$

The last inequality requires the assumption (8.5).

#### 8.4.9 Corrector for the clique estimator

With these results at hand, we can find a corrector for the clique estimator  $K$  for a randomly chosen vertex  $u$ .

**Theorem 8.5 (Corrector for clique estimator).** Let  $S$  be a multigraph supported on the vertex set  $F$ . Assume that the properties (8.4) and (8.5) hold. Draw  $u$  uniformly at random from  $F$ , and let  $K$  be the random estimator (8.7) for the clique induced by  $u$ . Then the random matrix  $\Phi(K - \mathbb{E}[K | u])$  has corrector

$$\frac{2f(\theta)}{|F|} \cdot I \quad \text{where} \quad f(\theta) = \exp\left(\frac{\theta^2/R}{1 - |\theta|/(3R)}\right) - 1.$$

The corrector is computed with respect to the randomness in the summands  $X_i$  and in the vertex  $u$ .

*Proof.* Proposition 8.4 gives a corrector of  $\Phi(K - \mathbb{E}[K | u])$  with respect to the randomness in the summands  $X_i$ . This corrector is

$$g(\theta) \Phi(\text{CLIQUE}(v, S)) \quad \text{where} \quad g(\theta) = \frac{\theta^2/(2R)}{1 - |\theta|/(3R)}.$$

We have shown that

$$\mathbb{E}_u \Phi(\text{CLIQUE}(u, S)) \leq \frac{4}{|F|} \cdot I \quad \text{and} \quad \|\Phi(\text{CLIQUE}(u, S))\| \leq 2.$$

Therefore, with respect to the random choice of the vertex  $u$ , the random matrix  $\Phi(\text{CLIQUE}(u, S))$  admits the corrector

$$\frac{4h(\theta)}{|F|} \cdot I \quad \text{where} \quad h(\theta) = \frac{e^{2\theta} - 1}{2}.$$

This is just the Chernoff corrector bound, Proposition 7.9. The result follows from the composition rule, Proposition 7.12, since  $f = 2(h \circ g)$ . ■

### 8.5 Analysis of SparseCholesky

We are finally prepared to establish that the SparseCholesky algorithm succeeds when we use the clique sampling procedure developed in Section 8.4. Fix a parameter  $\varepsilon \in (0, 1)$ . Our first goal is to prove that, with high probability,

$$\|\Phi(\mathbf{L}_n - \mathbf{L}_0)\| \leq \varepsilon.$$

Afterward, we must argue that the runtime of the algorithm is controlled.

#### 8.5.1 A stopping time

It is sufficient to show that, with high probability,

$$\max_{i=0,\dots,n} \|\Phi(\mathbf{L}_i - \mathbf{L}_0)\| \leq \varepsilon.$$

Let us define the stopping time

$$T = \min\{0 \leq i \leq n : \|\Phi(\mathbf{L}_i - \mathbf{L}_0)\| > \varepsilon\}.$$

If this event never occurs, then  $T = +\infty$ . For each  $i < T$ , observe that  $\|\Phi(\mathbf{L}_i)\| \leq 1 + \varepsilon$  because  $\|\Phi(\mathbf{L}_0)\| = 1$ .

We will consider the stopped martingale

$$\mathbf{Y}_i = \Phi(\mathbf{L}_{i \wedge T} - \mathbf{L}_0).$$

Clearly, it suffices to obtain a probability bound for the event that the stopped martingale exhibits a large deviation:

$$\max_{0 \leq i \leq n} \|\mathbf{Y}_i\| > \varepsilon.$$

We will treat the maximum and minimum eigenvalue parts of this spectral norm bound separately, but the arguments are symmetrical.

#### 8.5.2 The approximate Schur complements

The purpose of introducing the stopped martingale is to guarantee that the approximate Schur complements are uniformly bounded. Indeed, since  $\mathbf{0} \preceq \mathbf{S}_i \preceq \mathbf{L}_i$  for each  $i$ , we have the consequence that

$$\max_{i \leq T} \|\Phi(\mathbf{S}_{i-1})\| \leq \max_{i < T} \|\Phi(\mathbf{L}_i)\| \leq 2.$$

This condition delivers the uniform bound (8.5), irrespective of the choice of  $\varepsilon$ .

Moreover, the initial Laplacian  $\mathbf{S}_0$  consists of multiedges with leverage score bounded by  $1/R$ . At each step of the iteration, we remove some multiedges from  $\mathbf{S}_{i-1}$  and then add back a random clique estimator. Proposition 8.2 ensures that each multiedge in the clique estimator also has leverage score bounded by  $1/R$ . By induction, the assumption (8.4) holds in every iteration.



### 8.5.3 The corrector process

For  $i \geq 1$ , the difference sequence of the martingale  $\mathbf{Y}_i$  is

$$\mathbf{Y}_i - \mathbf{Y}_{i-1} = \begin{cases} \Phi(\mathbf{K}_i - \mathbb{E}_{i-1}[\mathbf{K}_i | \mathbf{u}_i]), & i \leq T \\ \mathbf{0}, & i > T. \end{cases}$$

For  $i \leq T$ , the matrix  $\mathbf{S}_{i-1}$  satisfies the conditions required to invoke Theorem 8.5. The support of  $\mathbf{S}_{i-1}$  has cardinality  $|\mathbf{F}_{i-1}| = n - i + 1$ . Therefore, the increment  $\mathbf{Y}_i - \mathbf{Y}_{i-1}$  has the corrector

$$\frac{2g(\theta)}{n - i + 1} \cdot \mathbf{I} \quad \text{where} \quad g(\theta) = \exp\left(\frac{\theta^2/R}{1 - |\theta|/(3R)}\right) - 1.$$

(This is computed conditional on all of the random choices up to step  $i - 1$ .) For  $i > T$ , we can take the corrector to be the zero matrix.

Therefore, owing to Proposition 7.10, the martingale admits the nonrandom corrector process

$$g \mathbf{W}_i = 2g \left[ \sum_{j=1}^{i \wedge T} \frac{1}{n - j + 1} \right] \cdot \mathbf{I} \leq 2g \log(en) \cdot \mathbf{I}.$$

To obtain the bound, we have summed the harmonic series up to  $j = n$ .

The same corrector is valid for the negation  $\{-\mathbf{Y}_i\}$  of the martingale, so we can obtain matching bounds for the maximum and minimum eigenvalues.

### 8.5.4 The martingale tail bound

Finally, we can bound the probability that the Laplacian martingale exhibits a large deviation. Set  $\sigma^2 = 2 \log(en)$ .

$$\begin{aligned} \mathbb{P} \{ \|\Phi(\mathbf{L}_n - \mathbf{L}_0)\| > \varepsilon \} &\leq \mathbb{P} \{ \exists i : \|\Phi(\mathbf{L}_i - \mathbf{L}_0)\| > \varepsilon \} \\ &= \mathbb{P} \{ \exists i : \|\mathbf{Y}_i\| > \varepsilon \} \\ &\leq \mathbb{P} \{ \exists i : \lambda_{\max}(\mathbf{Y}_i) \geq \varepsilon \} + \mathbb{P} \{ \exists i : \lambda_{\max}(-\mathbf{Y}_i) \geq \varepsilon \}. \end{aligned}$$

Indeed, the stopping time is triggered by the failure event, so we can pass to the stopped martingale. Then we split the spectral norm into eigenvalues so we can apply the master tail bound for matrix martingales.

We are in a good position to bound the last two probabilities.

$$\begin{aligned} \mathbb{P} \{ \exists i : \lambda_{\max}(\mathbf{Y}_i) \geq \varepsilon \} &= \mathbb{P} \{ \exists i : \lambda_{\max}(\mathbf{Y}_i) \geq \varepsilon \text{ and } \lambda_{\max}(\mathbf{W}_i) \leq \sigma^2 \} \\ &\leq n \cdot \inf_{\theta > 0} \exp(-\varepsilon\theta + g(\theta)\sigma^2). \end{aligned}$$

We have used the fact that  $\lambda_{\max}(\mathbf{W}_i) \leq \sigma^2$  always. The last inequality is a direct application of Theorem 7.4, the master tail bound for matrix martingales. Likewise,

$$\mathbb{P} \{ \exists i : \lambda_{\max}(-\mathbf{Y}_i) \geq \varepsilon \} \leq n \cdot \inf_{\theta > 0} \exp(-\varepsilon\theta + g(\theta)\sigma^2).$$

This bound also follows from Theorem 7.4.

Altogether, we determine that

$$\mathbb{P} \{ \|\Phi(\mathbf{L}_n - \mathbf{L}_0)\| > \varepsilon \} \leq 2n \cdot \inf_{\theta > 0} \exp(-\varepsilon\theta + 2g(\theta)\log(en)).$$

We may select the parameters

$$\theta = 2\varepsilon^{-1} \log(en) \quad \text{and} \quad R = \lceil 4\theta^2 \rceil = \lceil 16\varepsilon^{-2} \log^2(en) \rceil.$$

In this case,  $g(\theta) \leq 0.3$  for any  $n \geq 1$  and  $\varepsilon < 1$ . We have the overall bound

$$\mathbb{P} \{ \|\Phi(\mathbf{L}_n - \mathbf{L}_0)\| > \varepsilon \} \leq (en)^{-0.4}.$$

The probability bound can, of course, be improved by increasing the value of  $\theta$ . To do so, however, we must also increase the value of  $R$ , which means that we split the multiedges in the initial graph into more pieces.

### 8.5.5 The running time

Last, we must assess the running time of the **SparseCholesky** algorithm. The first step is to split the edges in the multigraph into  $R$  pieces to obtain a total of  $M = Rm$  multiedges. This step costs  $O(M)$  time and memory accesses.

At the outset, there are  $M = Rm$  multiedges in the graph. At each iteration, we eliminate a vertex by removing all the multiedges incident on that vertex and adding a clique with (at most) the same number of multiedges. As a result, the number of multiedges in the graph never increases above  $M$  at any iteration.

Now, in iteration  $i$ , we select a vertex  $u_i$  at random from the  $n - i + 1$  remaining vertices. In expectation, the number  $t_i$  of multiedges incident on  $u_i$  satisfies  $t_i \leq M/(n - i + 1)$ .

To sample the clique induced by  $u_i$ , we need to draw  $t_i$  samples from a probability mass on  $t_i$  points. This task can be accomplished in  $O(t_i)$  time overall [BP17]. The rest of the computation of the clique estimator and its introduction into the current Laplacian involve  $O(t_i)$  arithmetic and memory accesses.

Afterward, we remove the star induced by  $u_i$ , which also contains  $t_i$  multiedges. This operation involves  $O(t_i)$  arithmetic and memory accesses.

In summary, the expected running time of the algorithm is on the order of

$$M + \sum_{i=1}^n t_i \asymp Rm \sum_{i=1}^n \frac{1}{n - i + 1} \asymp Rm \log n.$$

To obtain an error of  $\varepsilon = 0.5$ , we can take  $R = \Theta(\log^2 n)$ . Therefore, the overall runtime is  $O(m \log^3(n))$ , in expectation.

### 8.5.6 The grand finale

To summarize, we have established the following result.

**Theorem 8.6 (Sparse Cholesky).** Let  $\mathbf{L} \in \mathbb{H}_V$  be the Laplacian matrix of a connected graph  $G$  on a set  $V$  of  $n$  vertices and with  $m$  weighted edges. The **SparseCholesky** algorithm produces a a morally lower-triangular matrix  $\mathbf{C} \in \mathbb{R}^{V \times V}$  that satisfies

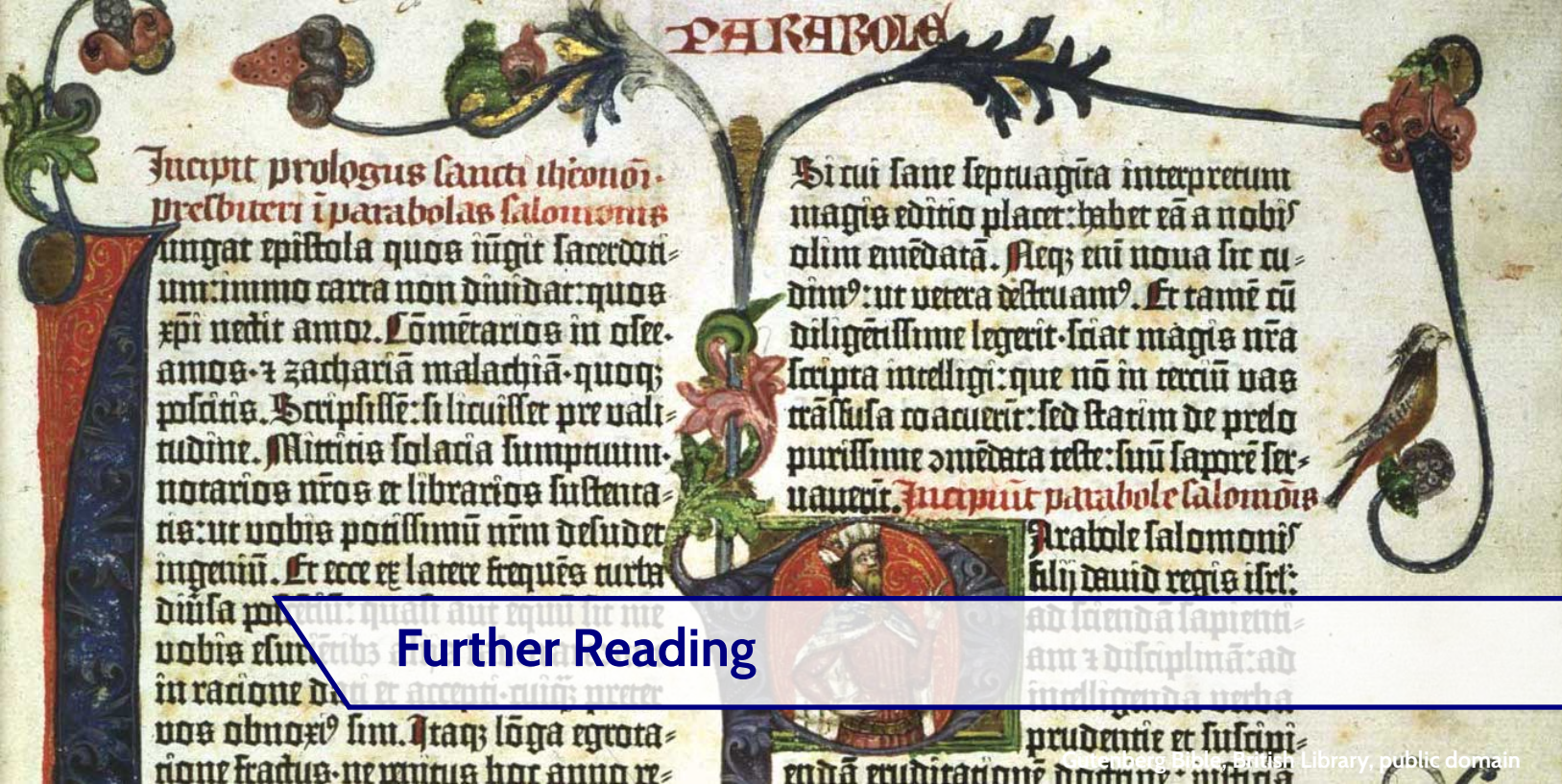
$$0.5 \mathbf{L} \preceq \mathbf{C} \mathbf{C}^* \preceq 1.5 \mathbf{L}.$$

The matrix  $\mathbf{C}$  has  $O(m \log^2 n)$  nonzero entries. The expected running time is  $O(m \log^3(n))$ .

In view of our discussion of preconditioned conjugate gradient, we arrive at an algorithmic approach for solving Laplacian linear systems.

**Corollary 8.7 (Laplacian systems).** Given the preconditioner  $\mathbf{L}$  computed by the Sparse-Cholesky algorithm, we can solve every consistent linear system in a graph Laplacian to relative error  $\varepsilon$  in the Dirichlet energy norm in time  $O(m \log^2(n) \log(1/\varepsilon))$ .

This is what we promised to prove. ■



Here is an incomplete collection of sources where you can learn more about contemporary random matrix theory and its applications.

#### Matrix concentration inequalities

The papers below develop results on how much a random matrix deviates from its mean in spectral norm. An important characteristic of these results is that they apply to a wide range of different types of random matrices, the constants are explicit (and reasonable), and the bounds are nonasymptotic.

- [Tro12] Joel A. Tropp, “User-friendly tail bounds for sums of random matrices.”  
This is a foundational paper that develops the modern approach to matrix concentration via the subadditivity of matrix cumulants. It contains a complete catalog of exponential inequalities for an independent sum of random matrices.
- [Tro11a] Joel A. Tropp, “Freedman’s inequality for matrix martingales.”  
This foundational paper is a follow-up to [Tro12] that develops the approach to matrix martingales using corrector processes. It was inspired by Roberto Oliveira’s paper [Oli10], which established a weaker version of the matrix Freedman inequality.
- [Tro15] Joel A. Tropp, “An introduction to matrix concentration inequalities.”  
My monograph gives a thorough introduction to matrix concentration for independent sums, including many applications, a complete proof of Lieb’s theorem, and an annotated bibliography of works on matrix concentration.
- [Tro16] Joel A. Tropp, “The expected norm of a sum of independent random matrices: An elementary approach.”

The matrix Rosenthal inequalities are moment inequalities for a sum of independent random matrices that strengthen the matrix Bernstein inequality. This paper proves the matrix Rosenthal inequalities using elementary arguments.

- [Mac+14] Lester Mackey et al., “Matrix concentration inequalities via the method of exchangeable pairs.”

This paper develops another approach to matrix concentration using Stein’s method. This approach is also more elementary than the approach using Lieb’s theorem, and it applies to some types of random matrices that are more general than independent sums or martingales.

- [PMT16] Daniel Paulin et al., “Efron–Stein inequalities for random matrices.”

This paper shows how to use Stein’s method to prove concentration inequalities for a matrix-valued function of independent random variables. This is potentially a very powerful result, but it has seen relatively few applications so far.

- [BH16] Afonso Bandeira and Ramon van Handel, “Sharp nonasymptotic bounds on the norm of random matrices with independent entries.”

This paper gives sharp bounds on the norm of a random matrix with independent entries, which is one of the most important random matrix models.

### Lower tail inequalities

The minimum eigenvalue of a sum of psd random matrices exhibits a totally different kind of behavior from the maximum eigenvalue. The following papers tackle this important problem.

- [KM15] Vladimir Koltchinskii and Shahar Mendelson, “Bounding the smallest singular value of a random matrix without concentration.”

This paper explains how to use the small-ball method to bound the smallest singular value of a random matrix with independent rows.

- [Tro16] Joel A. Tropp, “Convex recovery of a structured signal from independent random measurements.”

This expository work gives a simplified account of the small-ball method for controlling the minimum (conic) singular value of a random matrix with independent rows.

- [Oli16] Roberto Oliveira, “The lower tail of random quadratic forms with applications to ordinary least squares.”

This paper proves a lower tail inequality for a sum of independent, random psd matrices. The method is fascinating and strikingly different from other approaches.

### High-dimensional probability

Here are some surveys that I have found useful. They cover various aspects of high-dimensional probability with applications to modern random matrix theory.

- [FR13] Simon Foucart and Holger Rauhut, *A Mathematical Introduction to Compressive Sensing*.



This book develops a collection of methods for studying structured random matrices, and it proceeds from first principles.

- [Han16] Ramon van Handel, *Probability in High Dimensions*, 2016.  
These lecture notes contain a sophisticated mathematical treatment of high-dimensional probability, including some applications to random matrix theory.
- [Han17] Ramon van Handel, “Structured random matrices.”  
This accessible survey covers some very recent results about structured random matrices.
- [Ver18] Roman Vershynin, *High-Dimensional Probability*.  
This textbook gives an elementary introduction to the methods of high-dimensional probability, including some random matrix theory and many applications in data science.

#### Classical random matrix theory

Last, we mention a few resources for learning about the more established parts of random matrix theory.

- [Tao12] Terence Tao, *Topics in Random Matrix Theory*.  
This textbook gives an accessible introduction to the classical theory of random matrices.
- [Kem13] Todd Kemp, *Introduction to Random Matrix Theory*.  
These lecture notes provide another readable treatment of classical random matrix theory.
- [NSo6] Alexandru Nica and Roland Speicher, *Lectures on the Combinatorics of Free Probability*.  
This monograph introduces the theory of free probability from a combinatorial point of view.







## Bibliography

"Steacie Library," Wikimedia Commons

- [AW02] Rudolf Ahlswede and Andreas Winter. "Strong converse for identification via quantum channels". In: *IEEE Trans. Inform. Theory* 48.3 (2002), pages 569–579. DOI: [10.1109/18.985947](https://doi.org/10.1109/18.985947).
- [Axl15] Sheldon Axler. *Linear algebra done right*. Third. Undergraduate Texts in Mathematics. Springer, Cham, 2015, pages xviii+340. DOI: [10.1007/978-3-319-11080-6](https://doi.org/10.1007/978-3-319-11080-6).
- [BH16] Afonso S. Bandeira and Ramon van Handel. "Sharp nonasymptotic bounds on the norm of random matrices with independent entries". In: *Ann. Probab.* 44.4 (2016), pages 2479–2506. DOI: [10.1214/15-AOP1025](https://doi.org/10.1214/15-AOP1025).
- [BG13] Richard F. Bass and Karlheinz Gröchenig. "Relevant sampling of band-limited functions". In: *Illinois J. Math.* 57.1 (2013), pages 43–58. URL: <http://projecteuclid.org/euclid.ijm/1403534485>.
- [BSS14] Joshua Batson, Daniel A. Spielman, and Nikhil Srivastava. "Twice-Ramanujan sparsifiers". In: *SIAM Rev.* 56.2 (2014), pages 315–334. DOI: [10.1137/130949117](https://doi.org/10.1137/130949117).
- [Bha97] Rajendra Bhatia. *Matrix analysis*. Volume 169. Graduate Texts in Mathematics. Springer-Verlag, New York, 1997, pages xii+347. DOI: [10.1007/978-1-4612-0653-8](https://doi.org/10.1007/978-1-4612-0653-8). URL: <http://dx.doi.org/10.1007/978-1-4612-0653-8>.
- [Bha07] Rajendra Bhatia. *Positive definite matrices*. Princeton Series in Applied Mathematics. Princeton University Press, Princeton, NJ, 2007, pages x+254.

- [BLM13] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities*. A nonasymptotic theory of independence, With a foreword by Michel Ledoux. Oxford University Press, Oxford, 2013, pages x+481. DOI: [10.1093/acprof:oso/9780199535255.001.0001](https://doi.org/10.1093/acprof:oso/9780199535255.001.0001). URL: <http://dx.doi.org/10.1093/acprof:oso/9780199535255.001.0001>.
- [BP17] Karl Bringmann and Konstantinos Panagiotou. “Efficient sampling methods for discrete distributions”. In: *Algorithmica* 79.2 (2017), pages 484–508. DOI: [10.1007/s00453-016-0205-0](https://doi.org/10.1007/s00453-016-0205-0).
- [Buc01] Artur Buchholz. “Operator Khintchine inequality in non-commutative probability”. In: *Math. Ann.* 319.1 (2001), pages 1–16. DOI: [10.1007/PL00004425](https://doi.org/10.1007/PL00004425). URL: <http://dx.doi.org/10.1007/PL00004425>.
- [CC13] Xiaohong Chen and Timothy M. Christensen. “Optimal uniform convergence rates for sieve nonparametric instrumental variables regression”. Available at <http://arXiv.org/abs/1311.0412>. Nov. 2013.
- [Che+14] Yudong Chen et al. “Coherent matrix completion”. In: *Proc. 31st Intl. Conf. Machine Learning*. Beijing, 2014.
- [CSW12] Sin-Shuen Cheung, Anthony Man-Cho So, and Kuncheng Wang. “Linear matrix inequalities with stochastically dependent perturbations and applications to chance-constrained semidefinite optimization”. In: *SIAM J. Optim.* 22.4 (2012), pages 1394–1430. DOI: [10.1137/110822906](https://doi.org/10.1137/110822906). URL: <http://dx.doi.org/10.1137/110822906>.
- [CDL13] Albert Cohen, Mark A. Davenport, and Dany Leviatan. “On the stability and accuracy of least squares approximations”. In: *Found. Comput. Math.* 13.5 (2013), pages 819–834. DOI: [10.1007/s10208-013-9142-3](https://doi.org/10.1007/s10208-013-9142-3). URL: <http://dx.doi.org/10.1007/s10208-013-9142-3>.
- [CG14] Paul Constantine and David Gleich. “Computing active subspaces”. Available at <http://arXiv.org/abs/1408.0545>. Aug. 2014.
- [DKC13] Josip Djolonga, Andreas Krause, and Volkan Cevher. “High-Dimensional Gaussian Process Bandits”. In: *Advances in Neural Information Processing Systems* 26. Edited by C.J.C. Burges et al. Curran Associates, Inc., 2013, pages 1025–1033. URL: <http://papers.nips.cc/paper/5152-high-dimensional-gaussian-process-bandits.pdf>.
- [FSV12] Massimo Fornasier, Karin Schnass, and Jan Vybiral. “Learning functions of few arbitrary linear parameters in high dimensions”. In: *Found. Comput. Math.* 12.2 (2012), pages 229–262. DOI: [10.1007/s10208-012-9115-y](https://doi.org/10.1007/s10208-012-9115-y). URL: <http://dx.doi.org/10.1007/s10208-012-9115-y>.
- [FR13] Simon Foucart and Holger Rauhut. *A mathematical introduction to compressive sensing*. Applied and Numerical Harmonic Analysis. Birkhäuser/Springer, New York, 2013, pages xviii+625. DOI: [10.1007/978-0-8176-4948-7](https://doi.org/10.1007/978-0-8176-4948-7).
- [Fre75] David A. Freedman. “On tail probabilities for martingales”. In: *Ann. Probability* 3 (1975), pages 100–118. DOI: [10.1214/aop/1176996452](https://doi.org/10.1214/aop/1176996452).

- [GN51] Herman H. Goldstine and John von Neumann. “Numerical inverting of matrices of high order. II”. In: *Proc. Amer. Math. Soc.* 2 (1951), pages 188–202. DOI: [10.2307/2032484](https://doi.org/10.2307/2032484).
- [Grc11] Joseph F. Grcar. “John von Neumann’s analysis of Gaussian elimination and the origins of modern numerical analysis”. In: *SIAM Rev.* 53.4 (2011), pages 607–682. DOI: [10.1137/080734716](https://doi.org/10.1137/080734716).
- [Guh+18] Madeline Guha et al. “Fast state tomography with optimal error bounds”. Available at <http://arXiv.org/abs/1809.11162>. Sept. 2018.
- [Haa+17] Jidong Haah et al. “Sample-Optimal Tomography of Quantum States”. In: *IEEE Transactions on Information Theory* 63.9 (Sept. 2017), pages 5628–5641. DOI: [10.1109/TIT.2017.2719044](https://doi.org/10.1109/TIT.2017.2719044).
- [Han16] Ramon van Handel. “Probability in High Dimensions”. APC 550 Lecture Notes, Princeton University. Available at <https://web.math.princeton.edu/~rvan/APC550.pdf>. Dec. 2016.
- [Han17] Ramon van Handel. “Structured random matrices”. In: *Convexity and concentration*. Volume 161. IMA Vol. Math. Appl. Springer, New York, 2017, pages 107–156.
- [HO14] Nicholas J. A. Harvey and Neil Olver. “Pipage Rounding, Pessimistic Estimators and Matrix Concentration”. In: *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms*. SODA ’14. Portland, Oregon: SIAM, 2014, pages 926–945. URL: <http://dl.acm.org/citation.cfm?id=2634074.2634143>.
- [Higo2] Nicholas J. Higham. *Accuracy and stability of numerical algorithms*. Second. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002, pages xxx+680. DOI: [10.1137/1.9780898718027](https://doi.org/10.1137/1.9780898718027).
- [Hol12] Alexander S. Holevo. *Quantum systems, channels, information*. Volume 16. De Gruyter Studies in Mathematical Physics. A mathematical introduction. De Gruyter, Berlin, 2012, pages xiv+349. DOI: [10.1515/9783110273403](https://doi.org/10.1515/9783110273403). URL: <http://dx.doi.org/10.1515/9783110273403>.
- [HCG14] Qixing Huang, Yuxin Chen, and Leonidas Guibas. “Near-optimal joint object matching via convex relaxation”. In: *Proc. 31st Intl. Conf. Machine Learning*. Beijing, 2014.
- [KK12] Purushottam Kar and Harish Karnick. “Random Feature Maps for Dot Product Kernels”. In: *Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics*. Edited by Neil D. Lawrence and Mark Girolami. Volume 22. Proceedings of Machine Learning Research. La Palma, Canary Islands: PMLR, 21–23 Apr 2012, pages 583–591. URL: <http://proceedings.mlr.press/v22/kar12.html>.
- [Kem13] Todd Kemp. “Math 247A: Introduction to random matrix theory”. Available at <http://www.math.ucsd.edu/~tkemp/247A/247A.Notes.pdf>. 2013.



- [Kol11] Vladimir Koltchinskii. *Oracle inequalities in empirical risk minimization and sparse recovery problems*. Volume 2033. Lecture Notes in Mathematics. Lectures from the 38th Probability Summer School held in Saint-Flour, 2008, École d'Été de Probabilités de Saint-Flour. [Saint-Flour Probability Summer School]. Springer, Heidelberg, 2011, pages x+254. DOI: [10.1007/978-3-642-22147-7](https://doi.org/10.1007/978-3-642-22147-7). URL: <http://dx.doi.org/10.1007/978-3-642-22147-7>.
- [KM15] Vladimir Koltchinskii and Shahar Mendelson. “Bounding the smallest singular value of a random matrix without concentration”. In: *Int. Math. Res. Not. IMRN* 23 (2015), pages 12991–13008. DOI: [10.1093/imrn/rnv096](https://doi.org/10.1093/imrn/rnv096).
- [Kyn17] Rasmus Kyng. “Approximate Gaussian elimination”. PhD thesis. Yale University, 2017, page 120.
- [KS16] Rasmus Kyng and Sushant Sachdeva. “Approximate Gaussian elimination for Laplacians—fast, sparse, and simple”. In: *57th Annual IEEE Symposium on Foundations of Computer Science—FOCS 2016*. IEEE Computer Soc., Los Alamitos, CA, 2016, pages 573–582.
- [Kyn+16] Rasmus Kyng et al. “Sparsified Cholesky and multigrid solvers for connection Laplacians”. In: *STOC’16—Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing*. ACM, New York, 2016, pages 842–850.
- [Lie73] Elliott H. Lieb. “Convex trace functions and the Wigner-Yanase-Dyson conjecture”. In: *Advances in Math.* 11 (1973), pages 267–288. DOI: [10.1016/0001-8708\(73\)90011-X](https://doi.org/10.1016/0001-8708(73)90011-X).
- [Lop+14] David Lopez-Paz et al. “Randomized nonlinear component analysis”. In: *Proc. 31st Intl. Conf. Machine Learning*. Beijing, July 2014.
- [Lus86] Françoise Lust-Piquard. “Inégalités de Khintchine dans  $C_p$  ( $1 < p < \infty$ )”. In: *C. R. Acad. Sci. Paris Sér. I Math.* 303.7 (1986), pages 289–292.
- [LP91] Françoise Lust-Piquard and Gilles Pisier. “Noncommutative Khintchine and Paley inequalities”. In: *Ark. Mat.* 29.2 (1991), pages 241–260. DOI: [10.1007/BF02384340](https://doi.org/10.1007/BF02384340).
- [Mac+14] Lester Mackey et al. “Matrix concentration inequalities via the method of exchangeable pairs”. In: *Ann. Probab.* 42.3 (2014), pages 906–945. DOI: [10.1214/13-AOP892](https://doi.org/10.1214/13-AOP892). URL: <http://dx.doi.org/10.1214/13-AOP892>.
- [MB17] William B. March and George Biros. “Far-Field Compression for Fast Kernel Summation Methods in High Dimensions”. In: *Appl. Comput. Harmon. Anal.* 43.1 (July 2017), pages 39–75. DOI: [10.1016/j.acha.2015.09.007](https://doi.org/10.1016/j.acha.2015.09.007).
- [MKR12] Emilie Morvant, Sokol Koço, and Liva Ralaivola. “PAC-Bayesian Generalization Bound on Confusion Matrix for Multi-Class Classification”. In: *Proc. 29th Intl. Conf. Machine Learning*. Edinburgh, 2012.

- [NG47] John von Neumann and Herman H. Goldstine. “Numerical inverting of matrices of high order”. In: *Bull. Amer. Math. Soc.* 53 (1947), pages 1021–1099. DOI: [10.1090/S0002-9904-1947-08909-6](https://doi.org/10.1090/S0002-9904-1947-08909-6).
- [NSo6] Alexandru Nica and Roland Speicher. *Lectures on the combinatorics of free probability*. Volume 335. London Mathematical Society Lecture Note Series. Cambridge University Press, Cambridge, 2006, pages xvi+417. DOI: [10.1017/CB09780511735127](https://doi.org/10.1017/CB09780511735127).
- [Oli10] Roberto Imbuzeiro Oliveira. “The spectrum of random  $k$ -lifts of large graphs (with possibly large  $k$ )”. In: *J. Comb.* 1.3-4 (2010), pages 285–306. DOI: [10.4310/JOC.2010.v1.n3.a2](https://doi.org/10.4310/JOC.2010.v1.n3.a2). URL: <http://dx.doi.org/10.4310/JOC.2010.v1.n3.a2>.
- [Oli16] Roberto Imbuzeiro Oliveira. “The lower tail of random quadratic forms with applications to ordinary least squares”. In: *Probab. Theory Related Fields* 166.3-4 (2016), pages 1175–1194. DOI: [10.1007/s00440-016-0738-9](https://doi.org/10.1007/s00440-016-0738-9).
- [PMT16] Daniel Paulin, Lester Mackey, and Joel A. Tropp. “Efron-Stein inequalities for random matrices”. In: *Ann. Probab.* 44.5 (2016), pages 3431–3473. DOI: [10.1214/15-AOP1054](https://doi.org/10.1214/15-AOP1054). URL: <https://doi-org.clsproxy.library.caltech.edu/10.1214/15-AOP1054>.
- [PX97] Gilles Pisier and Quanhua Xu. “Non-commutative martingale inequalities”. In: *Comm. Math. Phys.* 189.3 (1997), pages 667–698. DOI: [10.1007/s002200050224](https://doi.org/10.1007/s002200050224).
- [RRo7] Ali Rahimi and Benjamin Recht. “Random features for large-scale kernel machines”. In: *Adv. Neural Information Processing Systems*. Vancouver, 2007.
- [Rud99] Mark Rudelson. “Random vectors in the isotropic position”. In: *J. Funct. Anal.* 164.1 (1999), pages 60–72. DOI: [10.1006/jfan.1998.3384](https://doi.org/10.1006/jfan.1998.3384). URL: <http://dx.doi.org/10.1006/jfan.1998.3384>.
- [SSo1] Bernhard Schölkopf and Alex Smola. *Learning with kernels*. Adaptive Computation and Machine Learning series. MIT Press, 2001.
- [Spi12] Daniel A. Spielman. “Spectral graph theory”. In: *Combinatorial scientific computing*. Chapman & Hall/CRC Comput. Sci. Ser. CRC Press, Boca Raton, FL, 2012, pages 495–524. DOI: [10.1201/b11644-19](https://doi.org/10.1201/b11644-19).
- [Spi] Daniel A. Spielman. *CPSC 662 / AMTH 561: Spectral Graph Theory*. URL: <http://www.cs.yale.edu/homes/spielman/561/syllabus.html> (visited on 06/29/2019).
- [SS11] Daniel A. Spielman and Nikhil Srivastava. “Graph sparsification by effective resistances”. In: *SIAM J. Comput.* 40.6 (2011), pages 1913–1926. DOI: [10.1137/080734029](https://doi.org/10.1137/080734029).
- [Tao12] Terence Tao. *Topics in random matrix theory*. Volume 132. Graduate Studies in Mathematics. American Mathematical Society, Providence, RI, 2012, pages x+282.

- [Tom74] Nicole Tomczak-Jaegermann. “The moduli of smoothness and convexity and the Rademacher averages of trace classes  $S_p(1 \leq p < \infty)$ ”. In: *Studia Math.* 50 (1974), pages 163–182.
- [TB97] Lloyd N. Trefethen and David Bau III. *Numerical linear algebra*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1997, pages xii+361. DOI: [10.1137/1.9780898719574](https://doi.org/10.1137/1.9780898719574).
- [Tro11a] Joel A. Tropp. “Freedman’s inequality for matrix martingales”. In: *Electron. Commun. Probab.* 16 (2011), pages 262–270. DOI: [10.1214/ECP.v16-1624](https://doi.org/10.1214/ECP.v16-1624).
- [Tro11b] Joel A. Tropp. “Improved analysis of the subsampled randomized Hadamard transform”. In: *Adv. Adapt. Data Anal.* 3.1-2 (2011), pages 115–126. DOI: [10.1142/S1793536911000787](https://doi.org/10.1142/S1793536911000787). URL: <http://dx.doi.org/10.1142/S1793536911000787>.
- [Tro12] Joel A. Tropp. “User-friendly tail bounds for sums of random matrices”. In: *Found. Comput. Math.* 12.4 (2012), pages 389–434. DOI: [10.1007/s10208-011-9099-z](https://doi.org/10.1007/s10208-011-9099-z).
- [Tro15] Joel A. Tropp. “An introduction to matrix concentration inequalities”. In: *Foundations and Trends in Machine Learning* 8.1–2 (May 2015), pages 1–230.
- [Tro16] Joel A. Tropp. “The expected norm of a sum of independent random matrices: an elementary approach”. In: *High dimensional probability VII*. Volume 71. Progr. Probab. Springer, [Cham], 2016, pages 173–202. DOI: [10.1007/978-3-319-40519-3\\_8](https://doi.org/10.1007/978-3-319-40519-3_8).
- [Ver18] Roman Vershynin. *High-dimensional probability*. Volume 47. Cambridge Series in Statistical and Probabilistic Mathematics. An introduction with applications in data science, With a foreword by Sara van de Geer. Cambridge University Press, Cambridge, 2018, pages xiv+284. DOI: [10.1017/9781108231596](https://doi.org/10.1017/9781108231596).
- [Wig] Yuval Wigderson. *Harmonic functions on graphs*. URL: <http://web.stanford.edu/~yuvalwig/math/teaching/HarmonicNotes.pdf> (visited on 06/29/2019).
- [Wil91] David Williams. *Probability with martingales*. Cambridge Mathematical Textbooks. Cambridge University Press, Cambridge, 1991, pages xvi+251. DOI: [10.1017/CB09780511813658](https://doi.org/10.1017/CB09780511813658).